

PCTWORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : C12N 15/12, C07K 14/47, C12Q 1/68, C07K 16/18		A2	(11) International Publication Number: WO 99/58675 (43) International Publication Date: 18 November 1999 (18.11.99)
(21) International Application Number: PCT/US99/10602			
(22) International Filing Date: 13 May 1999 (13.05.99)			
(30) Priority Data:			
60/085,426	14 May 1998 (14.05.98)	US	enue, San Francisco, CA 94116 (US). POT, David; 1565 5th Avenue #102, San Francisco, CA 94112 (US). KASSAM, Altaf; 2659 Harold Street, Oakland, CA 94602 (US). LAMSON, George; 232 Sandringham Drive, Moraga, CA 94556 (US). DRMANAC, Radoje; 850 East Greenwich Place, Palo Alto, CA 94303 (US). CRKVENJAKOV, Radomir, 762 Haverhill Drive, Sunnyvale, CA 94068 (US). DICKSON, Mark; 1411 Gabilan Drive #B, Hollister, CA 95025 (US). DRMANAC, Snezana; 850 East Greenwich Place, Palo Alto, CA 94303 (US). LABAT, Ivan; 140 Acalanes Drive, Sunnyvale, CA 94086 (US). LESHKOWITZ, Dena; 678 Durshire Way, Sunnyvale, CA 94087 (US). KITA, David; 899 Bounty Drive, Foster City, CA 94404 (US). GARCIA, Veronica; Apartment 412, 396 Año Nuevo, Sunnyvale, CA 94086 (US). JONES, Lee, William; 396 Año Nuevo #412, Sunnyvale, CA 94086 (US). STACHE-CRAIN, Birgit; 345 South Mary Avenue, Sunnyvale, CA 94086 (US).
60/085,537	15 May 1998 (15.05.98)	US	
60/085,696	15 May 1998 (15.05.98)	US	
60/105,234	21 October 1998 (21.10.98)	US	
60/105,877	27 October 1998 (27.10.98)	US	
(71) Applicants: CHIRON CORPORATION [US/US]; 4560 Horton Street – R440, Emeryville, CA 94608 (US). HYSEQ INC. [US/US]; 675 Almanor Avenue, Sunnyvale, CA 94086 (US).			
(72) Inventors: WILLIAMS, Lewis, T.; 3 Miroflores, Tiburon, CA 94920 (US). ESCOBEDO, Jaime; 1470 Lavona Road, Alamo, CA 94507 (US). INNIS, Michael, A.; 315 Constance Place, Moraga, CA 94556 (US). GARCIA, Pablo, Dominguez; 882 Chenery Street, San Francisco, CA 94131 (US). SUDDUTH-KLINGER, Julie; 280 Lexington Road, Kensington, CA 94707 (US). REINHARD, Christoph; 1633 Clinton Avenue, Alameda, CA 94501 (US). GIESKE, Klaus; Chausseestrasse 92, D-10115 Berlin (DE). RANDAZZO, Filippo; Apartment 403, 690 Chestnut Street, San Francisco, CA 94133 (US). KENNEDY, Giulia, C.; 360 Castenada Av-			
(74) Agent: BLACKBURN, Robert, P.; Chiron Corporation, P.O. Box 8097, Emeryville, CA 94662-8097 (US).			
(81) Designated States: AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, UZ, VN, YU, ZA, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).			
<p>Published <i>Without international search report and to be republished upon receipt of that report.</i></p>			
(54) Title: HUMAN GENES AND GENE EXPRESSION PRODUCTS V			
(57) Abstract			
<p>This invention relates to novel human polynucleotides and variants thereof, their encoded polypeptides and variants thereof, to genes corresponding to these polynucleotides and to proteins expressed by the genes. The invention also relates to diagnostic and therapeutic agents employing such novel human polynucleotides, their corresponding genes or gene products, e.g., these genes and proteins, including probes, antisense constructs, and antibodies.</p>			

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LJ	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		

HUMAN GENES AND GENE EXPRESSION PRODUCTS V

Field of the Invention

The present invention relates to polynucleotides of human origin and the encoded gene
5 products.

Background of the Invention

Identification of novel polynucleotides, particularly those that encode an expressed gene product, is important in the advancement of drug discovery, diagnostic technologies, and the understanding of the progression and nature of complex diseases such as cancer. Identification of
10 genes expressed in different cell types isolated from sources that differ in disease state or stage, developmental stage, exposure to various environmental factors, the tissue of origin, the species from which the tissue was isolated, and the like is key to identifying the genetic factors that are responsible for the phenotypes associated with these various differences.

This invention provides novel human polynucleotides, the polypeptides encoded by these
15 polynucleotides, and the genes and proteins corresponding to these novel polynucleotides.

Summary of the Invention

This invention relates to novel human polynucleotides and variants thereof, their encoded polypeptides and variants thereof, to genes corresponding to these polynucleotides and to proteins expressed by the genes. The invention also relates to diagnostics and therapeutics comprising such
20 novel human polynucleotides, their corresponding genes or gene products, including probes, antisense nucleotides, and antibodies. The polynucleotides of the invention correspond to a polynucleotide comprising the sequence information of at least one of SEQ ID NOS:1-2707.

Various aspects and embodiments of the invention will be readily apparent to the ordinarily skilled artisan upon reading the description provided herein.

25 Detailed Description of the Invention

The invention relates to polynucleotides comprising the disclosed nucleotide sequences, to full length cDNA, mRNA genomic sequences, and genes corresponding to these sequences and degenerate variants thereof, and to polypeptides encoded by the polynucleotides of the invention and polypeptide variants. The following detailed description describes the polynucleotide compositions
30 encompassed by the invention, methods for obtaining cDNA or genomic DNA encoding a full-length gene product, expression of these polynucleotides and genes, identification of structural motifs of the polynucleotides and genes, identification of the function of a gene product encoded by a gene corresponding to a polynucleotide of the invention, use of the provided polynucleotides as probes and in mapping and in tissue profiling, use of the corresponding polypeptides and other gene

products to raise antibodies, and use of the polynucleotides and their encoded gene products for therapeutic and diagnostic purposes.

Polynucleotide Compositions

The scope of the invention with respect to polynucleotide compositions includes, but is not necessarily limited to, polynucleotides having a sequence set forth in any one of SEQ ID NOS:1-2707; polynucleotides obtained from the biological materials described herein or other biological sources (particularly human sources) by hybridization under stringent conditions (particularly conditions of high stringency); genes corresponding to the provided polynucleotides; variants of the provided polynucleotides and their corresponding genes, particularly those variants that retain a biological activity of the encoded gene product (e.g., a biological activity ascribed to a gene product corresponding to the provided polynucleotides as a result of the assignment of the gene product to a protein family(ies) and/or identification of a functional domain present in the gene product). Other nucleic acid compositions contemplated by and within the scope of the present invention will be readily apparent to one of ordinary skill in the art when provided with the disclosure here.

“Polynucleotide” and “nucleic acid” as used herein with reference to nucleic acids of the composition is not intended to be limiting as to the length or structure of the nucleic acid unless specifically indicated.

The invention features polynucleotides that are expressed in human tissue, specifically human colon, breast, and/or lung tissue. Novel nucleic acid compositions of the invention of particular interest comprise a sequence set forth in any one of SEQ ID NOS:1-2707 or an identifying sequence thereof. An “identifying sequence” is a contiguous sequence of residues at least about 10 nt to about 20 nt in length, usually at least about 50 nt to about 100 nt in length, that uniquely identifies a polynucleotide sequence, e.g., exhibits less than 90%, usually less than about 80% to about 85% sequence identity to any contiguous nucleotide sequence of more than about 20 nt. Thus, the subject novel nucleic acid compositions include full length cDNAs or mRNAs that encompass an identifying sequence of contiguous nucleotides from any one of SEQ ID NOS: 1-2707.

The polynucleotides of the invention also include polynucleotides having sequence similarity or sequence identity. Nucleic acids having sequence similarity are detected by hybridization under low stringency conditions, for example, at 50°C and 10XSSC (0.9 M saline/0.09 M sodium citrate) and remain bound when subjected to washing at 55°C in 1XSSC. Sequence identity can be determined by hybridization under stringent conditions, for example, at 50°C or higher and 0.1XSSC (9 mM saline/0.9 mM sodium citrate). Hybridization methods and conditions are well known in the art, see, e.g., USPN 5,707,829. Nucleic acids that are substantially identical to the provided polynucleotide sequences, e.g. allelic variants, genetically altered versions of the gene,

etc. bind to the provided polynucleotide sequences (SEQ ID NOS:1-2707) under stringent hybridization conditions. By using probes, particularly labeled probes of DNA sequences, one can isolate homologous or related genes. The source of homologous genes can be any species, *e.g.* primate species, particularly human; rodents, such as rats and mice; canines, felines, bovines, 5 ovines, equines, yeast, nematodes, *etc.*

Preferably, hybridization is performed using at least 15 contiguous nucleotides (nt) of at least one of SEQ ID NOS:1-2707. That is, when at least 15 contiguous nt of one of the disclosed SEQ ID NOS. is used as a probe, the probe will preferentially hybridize with a nucleic acid comprising the complementary sequence, allowing the identification and retrieval of the nucleic acids that uniquely hybridize to the selected probe. Probes from more than one SEQ ID NO. can hybridize with the same nucleic acid if the cDNA from which they were derived corresponds to one mRNA. Probes of more than 15 nt can be used, *e.g.*, probes of from about 18 nt to about 100 nt, but 10 15 nt represents sufficient sequence for unique identification.

The polynucleotides of the invention also include naturally occurring variants of the 15 nucleotide sequences (*e.g.*, degenerate variants, allelic variants, *etc.*). Variants of the polynucleotides of the invention are identified by hybridization of putative variants with nucleotide sequences disclosed herein, preferably by hybridization under stringent conditions. For example, by using appropriate wash conditions, variants of the polynucleotides of the invention can be identified where the allelic variant exhibits at most about 25-30% base pair (bp) mismatches relative to the 20 selected polynucleotide probe. In general, allelic variants contain 15-25% bp mismatches, and can contain as little as even 5-15%, or 2-5%, or 1-2% bp mismatches, as well as a single bp mismatch.

The invention also encompasses homologs corresponding to the polynucleotides of SEQ ID 25 NOS:1-2707, where the source of homologous genes can be any mammalian species, *e.g.*, primate species, particularly human; rodents, such as rats; canines, felines, bovines, ovines, equines, yeast, nematodes, *etc.* Between mammalian species, *e.g.*, human and mouse, homologs generally have substantial sequence similarity, *e.g.*, at least 75% sequence identity, usually at least 90%, more usually at least 95% between nucleotide sequences. Sequence similarity is calculated based on a reference sequence, which may be a subset of a larger sequence, such as a conserved motif, coding region, flanking region, *etc.* A reference sequence will usually be at least about 18 contiguous nt 30 long, more usually at least about 30 nt long, and may extend to the complete sequence that is being compared. Algorithms for sequence analysis are known in the art, such as gapped BLAST, described in Altschul, et al. *Nucleic Acids Res.* (1997) 25:3389-3402.

In general, variants of the invention have a sequence identity greater than at least about 35 65%, preferably at least about 75%, more preferably at least about 85%, and can be greater than at least about 90% or more as determined by the Smith-Waterman homology search algorithm as

implemented in MPSRCH program (Oxford Molecular). For the purposes of this invention, a preferred method of calculating percent identity is the Smith-Waterman algorithm, using the following. Global DNA sequence identity must be greater than 65% as determined by the Smith-Waterman homology search algorithm as implemented in MPSRCH program (Oxford Molecular) 5 using an affine gap search with the following search parameters: gap open penalty, 12; and gap extension penalty, 1.

The subject nucleic acids can be cDNAs or genomic DNAs, as well as fragments thereof, particularly fragments that encode a biologically active gene product and/or are useful in the methods disclosed herein (e.g., in diagnosis, as a unique identifier of a differentially expressed gene 10 of interest, etc.). The term "cDNA" as used herein is intended to include all nucleic acids that share the arrangement of sequence elements found in native mature mRNA species, where sequence elements are exons and 3' and 5' non-coding regions. Normally mRNA species have contiguous exons, with the intervening introns, when present, being removed by nuclear RNA splicing, to create a continuous open reading frame encoding a polypeptide of the invention.

15 A genomic sequence of interest comprises the nucleic acid present between the initiation codon and the stop codon, as defined in the listed sequences, including all of the introns that are normally present in a native chromosome. It can further include the 3' and 5' untranslated regions found in the mature mRNA. It can further include specific transcriptional and translational regulatory sequences, such as promoters, enhancers, etc., including about 1 kb, but possibly more, of 20 flanking genomic DNA at either the 5' and 3' end of the transcribed region. The genomic DNA can be isolated as a fragment of 100 kbp or smaller, and substantially free of flanking chromosomal sequence. The genomic DNA flanking the coding region, either 3' and 5', or internal regulatory sequences as sometimes found in introns, contains sequences required for proper tissue, stage-specific, or disease-state specific expression.

25 The nucleic acid compositions of the subject invention can encode all or a part of the subject polypeptides. Double or single stranded fragments can be obtained from the DNA sequence by chemically synthesizing oligonucleotides in accordance with conventional methods, by restriction enzyme digestion, by PCR amplification, etc. Isolated polynucleotides and polynucleotide fragments of the invention comprise at least about 10, about 15, about 20, about 35, about 50, about 30 100, about 150 to about 200, about 250 to about 300, or about 350 contiguous nt selected from the polynucleotide sequences as shown in SEQ ID NOS:1-2707. For the most part, fragments will be of at least 15 nt, usually at least 18 nt or 25 nt, and up to at least about 50 contiguous nt in length or more. In a preferred embodiment, the polynucleotide molecules comprise a contiguous sequence of at least 12 nt selected from the group consisting of the polynucleotides shown in SEQ ID NOS:1- 35 2707.

Probes specific to the polynucleotides of the invention can be generated using the polynucleotide sequences disclosed in SEQ ID NOS:1-2707. The probes are preferably at least about a 12, 15, 16, 18, 20, 22, 24, or 25 nt fragment of a corresponding contiguous sequence of SEQ ID NOS:1-2707, and can be less than 2, 1, 0.5, 0.1, or 0.05 kb in length. The probes can be synthesized chemically or can be generated from longer polynucleotides using restriction enzymes. The probes can be labeled, for example, with a radioactive, biotinylated, or fluorescent tag. Preferably, probes are designed based upon an identifying sequence of a polynucleotide of one of SEQ ID NOS:1-2707. More preferably, probes are designed based on a contiguous sequence of one of the subject polynucleotides that remain unmasked following application of a masking program for masking low complexity (e.g., XBLAST) to the sequence.. *i.e.*.. one would select an unmasked region, as indicated by the polynucleotides outside the poly-n stretches of the masked sequence produced by the masking program.

The polynucleotides of the subject invention are isolated and obtained in substantial purity, generally as other than an intact chromosome. Usually, the polynucleotides, either as DNA or RNA, will be obtained substantially free of other naturally-occurring nucleic acid sequences, generally being at least about 50%, usually at least about 90% pure and are typically "recombinant", *e.g.*, flanked by one or more nucleotides with which it is not normally associated on a naturally occurring chromosome.

The polynucleotides of the invention can be provided as a linear molecule or within a circular molecule, and can be provided within autonomously replicating molecules (vectors) or within molecules without replication sequences. Expression of the polynucleotides can be regulated by their own or by other regulatory sequences known in the art. The polynucleotides of the invention can be introduced into suitable host cells using a variety of techniques available in the art, such as transferrin polycation-mediated DNA transfer, transfection with naked or encapsulated nucleic acids, liposome-mediated DNA transfer, intracellular transportation of DNA-coated latex beads, protoplast fusion, viral infection, electroporation, gene gun, calcium phosphate-mediated transfection, and the like.

The subject nucleic acid compositions can be used to, for example, produce polypeptides, as probes for the detection of mRNA of the invention in biological samples (*e.g.*, extracts of human cells) to generate additional copies of the polynucleotides, to generate ribozymes or antisense oligonucleotides, and as single stranded DNA probes or as triple-strand forming oligonucleotides. The probes described herein can be used to, for example, determine the presence or absence of the polynucleotide sequences as shown in SEQ ID NOS:1-2707 or variants thereof in a sample. These and other uses are described in more detail below.

Use of Polynucleotides to Obtain Full-Length cDNA, Gene, and Promoter Region

Full-length cDNA molecules comprising the disclosed polynucleotides are obtained as follows. A polynucleotide having a sequence of one of SEQ ID NOS:1-2707, or a portion thereof comprising at least 12, 15, 18, or 20 nt, is used as a hybridization probe to detect hybridizing members of a cDNA library using probe design methods, cloning methods, and clone selection techniques such as those described in USPN 5,654,173. Libraries of cDNA are made from selected tissues, such as normal or tumor tissue, or from tissues of a mammal treated with, for example, a pharmaceutical agent. Preferably, the tissue is the same as the tissue from which the polynucleotides of the invention were isolated, as both the polynucleotides described herein and the cDNA represent expressed genes. Most preferably, the cDNA library is made from the biological material described herein in the Examples. The choice of cell type for library construction can be made after the identity of the protein encoded by the gene corresponding to the polynucleotide of the invention is known. This will indicate which tissue and cell types are likely to express the related gene, and thus represent a suitable source for the mRNA for generating the cDNA. Where the provided polynucleotides are isolated from cDNA libraries, the libraries are prepared from mRNA of human colon cells, more preferably, human colon cancer cells, even more preferably, from a highly metastatic colon cell, Km12L4-A.

Techniques for producing and probing nucleic acid sequence libraries are described, for example, in Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, 2nd Ed., (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY. The cDNA can be prepared by using primers based on sequence from SEQ ID NOS:1-2707. In one embodiment, the cDNA library can be made from only poly-adenylated mRNA. Thus, poly-T primers can be used to prepare cDNA from the mRNA.

Members of the library that are larger than the provided polynucleotides, and preferably that encompass the complete coding sequence of the native message, are obtained. In order to confirm that the entire cDNA has been obtained, RNA protection experiments are performed as follows. Hybridization of a full-length cDNA to an mRNA will protect the RNA from RNase degradation. If the cDNA is not full length, then the portions of the mRNA that are not hybridized will be subject to RNase degradation. This is assayed, as is known in the art, by changes in electrophoretic mobility on polyacrylamide gels, or by detection of released monoribonucleotides. Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, 2nd Ed., (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY. In order to obtain additional sequences 5' to the end of a partial cDNA, 5' RACE (*PCR Protocols: A Guide to Methods and Applications*, (1990) Academic Press, Inc.) can be performed.

Genomic DNA is isolated using the provided polynucleotides in a manner similar to the isolation of full-length cDNAs. Briefly, the provided polynucleotides, or portions thereof, are used as probes to libraries of genomic DNA. Preferably, the library is obtained from the cell type that

was used to generate the polynucleotides of the invention, but this is not essential. Most preferably, the genomic DNA is obtained from the biological material described herein in the Examples. Such libraries can be in vectors suitable for carrying large segments of a genome, such as P1 or YAC, as described in detail in Sambrook *et al.*, 9.4-9.30. In addition, genomic sequences can be isolated 5 from human BAC libraries, which are commercially available from Research Genetics, Inc., Huntsville, Alabama, USA, for example. In order to obtain additional 5' or 3' sequences, chromosome walking is performed, as described in Sambrook *et al.*, such that adjacent and overlapping fragments of genomic DNA are isolated. These are mapped and pieced together, as is known in the art, using restriction digestion enzymes and DNA ligase.

10 Using the polynucleotide sequences of the invention, corresponding full-length genes can be isolated using both classical and PCR methods to construct and probe cDNA libraries. Using either method, Northern blots, preferably, are performed on a number of cell types to determine which cell lines express the gene of interest at the highest level. Classical methods of constructing cDNA libraries are taught in Sambrook *et al.*, *supra*. With these methods, cDNA can be produced from 15 mRNA and inserted into viral or expression vectors. Typically, libraries of mRNA comprising poly(A) tails can be produced with poly(T) primers. Similarly, cDNA libraries can be produced using the instant sequences as primers.

PCR methods are used to amplify the members of a cDNA library that comprise the desired insert. In this case, the desired insert will contain sequence from the full length cDNA that 20 corresponds to the instant polynucleotides. Such PCR methods include gene trapping and RACE methods. Gene trapping entails inserting a member of a cDNA library into a vector. The vector then is denatured to produce single stranded molecules. Next, a substrate-bound probe, such a biotinylated oligo, is used to trap cDNA inserts of interest. Biotinylated probes can be linked to an avidin-bound solid substrate. PCR methods can be used to amplify the trapped cDNA. To trap 25 sequences corresponding to the full length genes, the labeled probe sequence is based on the polynucleotide sequences of the invention. Random primers or primers specific to the library vector can be used to amplify the trapped cDNA. Such gene trapping techniques are described in Gruber *et al.*, WO 95/04745 and Gruber *et al.*, USPN 5,500,356. Kits are commercially available to perform gene trapping experiments from, for example, Life Technologies, Gaithersburg, Maryland, USA.

30 "Rapid amplification of cDNA ends," or RACE, is a PCR method of amplifying cDNAs from a number of different RNAs. The cDNAs are ligated to an oligonucleotide linker, and amplified by PCR using two primers. One primer is based on sequence from the instant polynucleotides, for which full length sequence is desired, and a second primer comprises sequence that hybridizes to the oligonucleotide linker to amplify the cDNA. A description of this methods is 35 reported in WO 97/19110. In preferred embodiments of RACE, a common primer is designed to

anneal to an arbitrary adaptor sequence ligated to cDNA ends (Apte and Siebert, *Biotechniques* (1993) 15:890-893; Edwards *et al.*, *Nuc. Acids Res.* (1991) 19:5227-5232). When a single gene-specific RACE primer is paired with the common primer, preferential amplification of sequences between the single gene specific primer and the common primer occurs. Commercial cDNA pools modified for use in RACE are available.

Another PCR-based method generates full-length cDNA library with anchored ends without needing specific knowledge of the cDNA sequence. The method uses lock-docking primers (I-VI), where one primer, poly TV (I-III) locks over the polyA tail of eukaryotic mRNA producing first strand synthesis and a second primer, polyGH (IV-VI) locks onto the polyC tail added by terminal 10 deoxynucleotidyl transferase (TdT) (see, e.g., WO 96/40998).

The promoter region of a gene generally is located 5' to the initiation site for RNA polymerase II. Hundreds of promoter regions contain the "TATA" box, a sequence such as TATTA or TATAA, which is sensitive to mutations. The promoter region can be obtained by performing 5' RACE using a primer from the coding region of the gene. Alternatively, the cDNA can be used as a 15 probe for the genomic sequence, and the region 5' to the coding region is identified by "walking up." If the gene is highly expressed or differentially expressed, the promoter from the gene can be of use in a regulatory construct for a heterologous gene.

Once the full-length cDNA or gene is obtained, DNA encoding variants can be prepared by site-directed mutagenesis, described in detail in Sambrook *et al.*, 15.3-15.63. The choice of codon or 20 nucleotide to be replaced can be based on disclosure herein on optional changes in amino acids to achieve altered protein structure and/or function.

As an alternative method to obtaining DNA or RNA from a biological material, nucleic acid comprising nucleotides having the sequence of one or more polynucleotides of the invention can be synthesized. Thus, the invention encompasses nucleic acid molecules ranging in length from 15 nt 25 (corresponding to at least 15 contiguous nt of one of SEQ ID NOS:1-2707) up to a maximum length suitable for one or more biological manipulations, including replication and expression, of the nucleic acid molecule. The invention includes but is not limited to (a) nucleic acid having the size of a full gene, and comprising at least one of SEQ ID NOS:1-2707; (b) the nucleic acid of (a) also comprising at least one additional gene, operably linked to permit expression of a fusion protein; (c) 30 an expression vector comprising (a) or (b); (d) a plasmid comprising (a) or (b); and (e) a recombinant viral particle comprising (a) or (b). Once provided with the polynucleotides disclosed herein, construction or preparation of (a) - (e) are well within the skill in the art.

The sequence of a nucleic acid comprising at least 15 contiguous nt of at least any one of SEQ ID NOS:1-2707, preferably the entire sequence of at least any one of SEQ ID NOS:1-2707, is 35 not limited and can be any sequence of A, T, G, and/or C (for DNA) and A, U, G, and/or C (for

RNA) or modified bases thereof, including inosine and pseudouridine. The choice of sequence will depend on the desired function and can be dictated by coding regions desired, the intron-like regions desired, and the regulatory regions desired. Where the entire sequence of any one of SEQ ID NOS:1-2707 is within the nucleic acid, the nucleic acid obtained is referred to herein as a 5 polynucleotide comprising the sequence of any one of SEQ ID NOS:1-2707.

Expression of Polypeptide Encoded by Full-Length cDNA or Full-Length Gene

The provided polynucleotides (e.g., a polynucleotide having a sequence of one of SEQ ID NOS:1-2707), the corresponding cDNA, or the full-length gene is used to express a partial or complete gene product. Constructs of polynucleotides having sequences of SEQ ID NOS:1-2707 10 can also be generated synthetically. Alternatively, single-step assembly of a gene and entire plasmid from large numbers of oligodeoxyribonucleotides is described by, e.g., Stemmer *et al.*, *Gene (Amsterdam)* (1995) 164(1):49-53. In this method, assembly PCR (the synthesis of long DNA sequences from large numbers of oligodeoxyribonucleotides (oligos)) is described. The method is derived from DNA shuffling (Stemmer, *Nature* (1994) 370:389-391), and does not rely on DNA 15 ligase, but instead relies on DNA polymerase to build increasingly longer DNA fragments during the assembly process.

Appropriate polynucleotide constructs are purified using standard recombinant DNA techniques as described in, for example, Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, 2nd Ed., (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY, and under current regulations 20 described in United States Dept. of HHS, National Institute of Health (NIH) Guidelines for Recombinant DNA Research. The gene product encoded by a polynucleotide of the invention is expressed in any expression system, including, for example, bacterial, yeast, insect, amphibian and mammalian systems. Vectors, host cells and methods for obtaining expression in same are well known in the art. Suitable vectors and host cells are described in USPN 5,654,173.

25 Polynucleotide molecules comprising a polynucleotide sequence provided herein are generally propagated by placing the molecule in a vector. Viral and non-viral vectors are used, including plasmids. The choice of plasmid will depend on the type of cell in which propagation is desired and the purpose of propagation. Certain vectors are useful for amplifying and making large amounts of the desired DNA sequence. Other vectors are suitable for expression in cells in culture. 30 Still other vectors are suitable for transfer and expression in cells in a whole animal or person. The choice of appropriate vector is well within the skill of the art. Many such vectors are available commercially. Methods for preparation of vectors comprising a desired sequence are well known in the art.

35 The polynucleotides set forth in SEQ ID NOS:1-2707 or their corresponding full-length polynucleotides are linked to regulatory sequences as appropriate to obtain the desired expression

properties. These can include promoters (attached either at the 5' end of the sense strand or at the 3' end of the antisense strand), enhancers, terminators, operators, repressors, and inducers. The promoters can be regulated or constitutive. In some situations it may be desirable to use conditionally active promoters, such as tissue-specific or developmental stage-specific promoters.

5 These are linked to the desired nucleotide sequence using the techniques described above for linkage to vectors. Any techniques known in the art can be used.

When any of the above host cells, or other appropriate host cells or organisms, are used to replicate and/or express the polynucleotides or nucleic acids of the invention, the resulting replicated nucleic acid, RNA, expressed protein or polypeptide, is within the scope of the invention as a 10 product of the host cell or organism. The product is recovered by any appropriate means known in the art.

Once the gene corresponding to a selected polynucleotide is identified, its expression can be regulated in the cell to which the gene is native. For example, an endogenous gene of a cell can be regulated by an exogenous regulatory sequence as disclosed in USPN 5,641,670.

15

Identification of Functional and Structural Motifs of Novel Genes Screening Against Publicly Available Databases

Translations of the nucleotide sequence of the provided polynucleotides, cDNAs or full genes can be aligned with individual known sequences. Similarity with individual sequences can be 20 used to determine the activity of the polypeptides encoded by the polynucleotides of the invention. Also, sequences exhibiting similarity with more than one individual sequence can exhibit activities that are characteristic of either or both individual sequences.

The full length sequences and fragments of the polynucleotide sequences of the nearest neighbors can be used as probes and primers to identify and isolate the full length sequence 25 corresponding to provided polynucleotides. The nearest neighbors can indicate a tissue or cell type to be used to construct a library for the full-length sequences corresponding to the provided polynucleotides.

Typically, a selected polynucleotide is translated in all six frames to determine the best alignment with the individual sequences. The sequences disclosed herein in the Sequence Listing 30 are in a 5' to 3' orientation and translation in three frames can be sufficient (with a few specific exceptions as described in the Examples). These amino acid sequences are referred to, generally, as query sequences, which will be aligned with the individual sequences. Databases with individual sequences are described in "Computer Methods for Macromolecular Sequence Analysis" *Methods in Enzymology* (1996) 266, Doolittle, Academic Press, Inc., a division of Harcourt Brace & Co., San 35 Diego, California, USA. Databases include GenBank, EMBL, and DNA Database of Japan (DDBJ).

Query and individual sequences can be aligned using the methods and computer programs described above, and include BLAST 2.0, available over the world wide web at <http://www.ncbi.nlm.nih.gov/BLAST/>. See also Altschul, et al. *Nucleic Acids Res.* (1997) 25:3389-3402. Another alignment algorithm is Fasta, available in the Genetics Computing Group (GCG) package, Madison, Wisconsin, USA, a wholly owned subsidiary of Oxford Molecular Group, Inc. Other techniques for alignment are described in Doolittle, *supra*. Preferably, an alignment program that permits gaps in the sequence is utilized to align the sequences. The Smith-Waterman is one type of algorithm that permits gaps in sequence alignments. See *Meth. Mol. Biol.* (1997) 70: 173-187. Also, the GAP program using the Needleman and Wunsch alignment method can be utilized to align sequences. An alternative search strategy uses MPSRCH software, which runs on a MASPAR computer. MPSRCH uses a Smith-Waterman algorithm to score sequences on a massively parallel computer. This approach improves ability to identify sequences that are distantly related matches, and is especially tolerant of small gaps and nucleotide sequence errors. Amino acid sequences encoded by the provided polynucleotides can be used to search both protein and DNA databases.

Incorporated herein by reference are all sequences that have been made public as of the filing date of this application by any of the DNA or protein sequence databases, including the patent databases (e.g., GeneSeq). Also incorporated by reference are those sequences that have been submitted to these databases as of the filing date of the present application but not made public until after the filing date of the present application.

Results of individual and query sequence alignments can be divided into three categories: high similarity, weak similarity, and no similarity. Individual alignment results ranging from high similarity to weak similarity provide a basis for determining polypeptide activity and/or structure. Parameters for categorizing individual results include: percentage of the alignment region length where the strongest alignment is found, percent sequence identity, and p value. The percentage of the alignment region length is calculated by counting the number of residues of the individual sequence found in the region of strongest alignment, e.g., contiguous region of the individual sequence that contains the greatest number of residues that are identical to the residues of the corresponding region of the aligned query sequence. This number is divided by the total residue length of the query sequence to calculate a percentage. For example, a query sequence of 20 amino acid residues might be aligned with a 20 amino acid region of an individual sequence. The individual sequence might be identical to amino acid residues 5, 9-15, and 17-19 of the query sequence. The region of strongest alignment is thus the region stretching from residue 9-19, an 11 amino acid stretch. The percentage of the alignment region length is: 11 (length of the region of strongest alignment) divided by (query sequence length) 20 or 55%.

Percent sequence identity is calculated by counting the number of amino acid matches between the query and individual sequence and dividing total number of matches by the number of residues of the individual sequences found in the region of strongest alignment. Thus, the percent identity in the example above would be 10 matches divided by 11 amino acids, or approximately.

5 90.9%

P value is the probability that the alignment was produced by chance. For a single alignment, the p value can be calculated according to Karlin *et al.*, *Proc. Natl. Acad. Sci.* (1990) 87:2264 and Karlin *et al.*, *Proc. Natl. Acad. Sci.* (1993) 90. The p value of multiple alignments using the same query sequence can be calculated using an heuristic approach described in Altschul *et al.*, *Nat. Genet.* (1994) 6:119. Alignment programs such as BLAST program can calculate the p value. See also Altschul *et al.*, *Nucleic Acids Res.* (1997) 25:3389-3402.

Another factor to consider for determining identity or similarity is the location of the similarity or identity. Strong local alignment can indicate similarity even if the length of alignment is short. Sequence identity scattered throughout the length of the query sequence also can indicate a 15 similarity between the query and profile sequences. The boundaries of the region where the sequences align can be determined according to Doolittle, *supra*; BLAST 2.0 (see, e.g., Altschul, *et al.* *Nucleic Acids Res.* (1997) 25:3389-3402) or FAST programs; or by determining the area where sequence identity is highest.

High Similarity. In general, in alignment results considered to be of high similarity, the 20 percent of the alignment region length is typically at least about 55% of total length query sequence; more typically, at least about 58%; even more typically; at least about 60% of the total residue length of the query sequence. Usually, percent length of the alignment region can be as much as about 62%; more usually, as much as about 64%; even more usually, as much as about 66%. Further, for high similarity, the region of alignment, typically, exhibits at least about 75% of 25 sequence identity; more typically, at least about 78%; even more typically; at least about 80% sequence identity. Usually, percent sequence identity can be as much as about 82%; more usually, as much as about 84%; even more usually, as much as about 86%.

The p value is used in conjunction with these methods. If high similarity is found, the query sequence is considered to have high similarity with a profile sequence when the p value is less than 30 or equal to about 10^{-2} ; more usually; less than or equal to about 10^{-3} ; even more usually; less than or equal to about 10^{-4} . More typically, the p value is no more than about 10^{-5} ; more typically; no more than or equal to about 10^{-10} ; even more typically; no more than or equal to about 10^{-15} for the query sequence to be considered high similarity.

Weak Similarity. In general, where alignment results considered to be of weak similarity, there is no minimum percent length of the alignment region nor minimum length of alignment. A better showing of weak similarity is considered when the region of alignment is, typically, at least about 15 amino acid residues in length; more typically, at least about 20; even more typically: at 5 least about 25 amino acid residues in length. Usually, length of the alignment region can be as much as about 30 amino acid residues; more usually, as much as about 40; even more usually, as much as about 60 amino acid residues. Further, for weak similarity, the region of alignment, typically, exhibits at least about 35% of sequence identity; more typically, at least about 40%; even more typically: at least about 45% sequence identity. Usually, percent sequence identity can be as much 10 as about 50%; more usually, as much as about 55%; even more usually, as much as about 60%.

If low similarity is found, the query sequence is considered to have weak similarity with a profile sequence when the p value is usually less than or equal to about 10^{-2} ; more usually: less than or equal to about 10^{-3} ; even more usually; less than or equal to about 10^{-4} . More typically, the p value is no more than about 10^{-5} ; more usually; no more than or equal to about 10^{-10} ; even more 15 usually; no more than or equal to about 10^{-15} for the query sequence to be considered weak similarity.

Similarity Determined by Sequence Identity Alone. Sequence identity alone can be used to determine similarity of a query sequence to an individual sequence and can indicate the activity of the sequence. Such an alignment, preferably, permits gaps to align sequences. Typically, the query 20 sequence is related to the profile sequence if the sequence identity over the entire query sequence is at least about 15%; more typically, at least about 20%; even more typically, at least about 25%; even more typically, at least about 50%. Sequence identity alone as a measure of similarity is most useful when the query sequence is usually, at least 80 residues in length; more usually, 90 residues; even more usually, at least 95 amino acid residues in length. More typically, similarity can be concluded 25 based on sequence identity alone when the query sequence is preferably 100 residues in length: more preferably, 120 residues in length; even more preferably, 150 amino acid residues in length.

Alignments with Profile and Multiple Aligned Sequences. Translations of the provided polynucleotides can be aligned with amino acid profiles that define either protein families or common motifs. Also, translations of the provided polynucleotides can be aligned to multiple 30 sequence alignments (MSA) comprising the polypeptide sequences of members of protein families or motifs. Similarity or identity with profile sequences or MSAs can be used to determine the activity of the gene products (e.g., polypeptides) encoded by the provided polynucleotides or

corresponding cDNA or genes. For example, sequences that show an identity or similarity with a chemokine profile or MSA can exhibit chemokine activities.

Profiles can be designed manually by (1) creating an MSA, which is an alignment of the amino acid sequence of members that belong to the family and (2) constructing a statistical representation of the alignment. Such methods are described, for example, in Birney *et al.*, *Nucl. Acid Res.* (1996) 24(14): 2730-2739. MSAs of some protein families and motifs are publicly available. For example, <http://genome.wustl.edu/Pfam/> includes MSAs of 547 different families and motifs. These MSAs are described also in Sonnhammer *et al.*, *Proteins* (1997) 28: 405-420. Other sources over the world wide web include the site at <http://www.embl-heidelberg.de/argos/ali/ali.html>; alternatively, a message can be sent to ALI@EMBL-HEIDELBERG.DE for the information. A brief description of these MSAs is reported in Pascarella *et al.*, *Prot. Eng.* (1996) 9(3):249-251. Techniques for building profiles from MSAs are described in Sonnhammer *et al.*, *supra*; Birney *et al.*, *supra*; and "Computer Methods for Macromolecular Sequence Analysis," *Methods in Enzymology* (1996) 266. Doolittle, Academic Press, Inc., San Diego, California, USA.

Similarity between a query sequence and a protein family or motif can be determined by (a) comparing the query sequence against the profile and/or (b) aligning the query sequence with the members of the family or motif. Typically, a program such as Searchwise is used to compare the query sequence to the statistical representation of the multiple alignment, also known as a profile (see Birney *et al.*, *supra*). Other techniques to compare the sequence and profile are described in Sonnhammer *et al.*, *supra* and Doolittle, *supra*.

Next, methods described by Feng *et al.*, *J. Mol. Evol.* (1987) 25:351 and Higgins *et al.*, *CABIOS* (1989) 5:151 can be used to align the query sequence with the members of a family or motif, also known as a MSA. Sequence alignments can be generated using any of a variety of software tools. Examples include PileUp, which creates a multiple sequence alignment, and is described in Feng *et al.*, *J. Mol. Evol.* (1987) 25:351. Another method, GAP, uses the alignment method of Needleman *et al.*, *J. Mol. Biol.* (1970) 48:443. GAP is best suited for global alignment of sequences. A third method, BestFit, functions by inserting gaps to maximize the number of matches using the local homology algorithm of Smith *et al.*, *Adv. Appl. Math.* (1981) 2:482. In general, the following factors are used to determine if a similarity between a query sequence and a profile or MSA exists:

(1) number of conserved residues found in the query sequence, (2) percentage of conserved residues found in the query sequence, (3) number of frameshifts, and (4) spacing between conserved residues.

Some alignment programs that both translate and align sequences can make any number of frameshifts when translating the nucleotide sequence to produce the best alignment. The fewer frameshifts needed to produce an alignment, the stronger the similarity or identity between the query and profile or MSAs. For example, a weak similarity resulting from no frameshifts can be a better

indication of activity or structure of a query sequence, than a strong similarity resulting from two frameshifts. Preferably, three or fewer frameshifts are found in an alignment; more preferably two or fewer frameshifts; even more preferably, one or fewer frameshifts; even more preferably, no frameshifts are found in an alignment of query and profile or MSAs.

5 Conserved residues are those amino acids found at a particular position in all or some of the family or motif members. Alternatively, a position is considered conserved if only a certain class of amino acids is found in a particular position in all or some of the family members. For example, the N-terminal position can contain a positively charged amino acid, such as lysine, arginine, or histidine.

10 Typically, a residue of a polypeptide is conserved when a class of amino acids or a single amino acid is found at a particular position in at least about 40% of all class members; more typically, at least about 50%; even more typically, at least about 60% of the members. Usually, a residue is conserved when a class or single amino acid is found in at least about 70% of the members of a family or motif; more usually, at least about 80%; even more usually, at least about 90%; even 15 more usually, at least about 95%.

20 A residue is considered conserved when three unrelated amino acids are found at a particular position in the some or all of the members; more usually, two unrelated amino acids. These residues are conserved when the unrelated amino acids are found at particular positions in at least about 40% of all class member; more typically, at least about 50%; even more typically, at least about 60% of the members. Usually, a residue is conserved when a class or single amino acid is found in at least about 70% of the members of a family or motif; more usually, at least about 80%; even more usually, at least about 90%; even more 25 usually, at least about 95%.

30 A query sequence has similarity to a profile or MSA when the query sequence comprises at least about 25% of the conserved residues of the profile or MSA; more usually, at least about 30%; even more usually; at least about 40%. Typically, the query sequence has a stronger similarity to a profile sequence or MSA when the query sequence comprises at least about 45% of the conserved residues of the profile or MSA; more typically, at least about 50%; even more typically; at least about 55%.

Identification of Secreted & Membrane-Bound Polypeptides

35 Both secreted and membrane-bound polypeptides of the present invention are of particular interest. For example, levels of secreted polypeptides can be assayed in body fluids that are convenient, such as blood, plasma, serum, and other body fluids such as urine, prostatic fluid and semen. Membrane-bound polypeptides are useful for constructing vaccine antigens or inducing an immune response. Such antigens would comprise all or part of the extracellular region of the membrane-bound polypeptides. Because both secreted and membrane-bound polypeptides comprise

a fragment of contiguous hydrophobic amino acids. hydrophobicity predicting algorithms can be used to identify such polypeptides.

A signal sequence is usually encoded by both secreted and membrane-bound polypeptide genes to direct a polypeptide to the surface of the cell. The signal sequence usually comprises a stretch of hydrophobic residues. Such signal sequences can fold into helical structures. Membrane-bound polypeptides typically comprise at least one transmembrane region that possesses a stretch of hydrophobic amino acids that can transverse the membrane. Some transmembrane regions also exhibit a helical structure. Hydrophobic fragments within a polypeptide can be identified by using computer algorithms. Such algorithms include Hopp & Woods, *Proc. Natl. Acad. Sci. USA* (1981) 78:3824-3828; Kyte & Doolittle, *J. Mol. Biol.* (1982) 157: 105-132; and RAOAR algorithm. Degli Esposti *et al.*, *Eur. J. Biochem.* (1990) 190: 207-219.

Another method of identifying secreted and membrane-bound polypeptides is to translate the polynucleotides of the invention in all six frames and determine if at least 8 contiguous hydrophobic amino acids are present. Those translated polypeptides with at least 8; more typically, 10: even more typically, 12 contiguous hydrophobic amino acids are considered to be either a putative secreted or membrane bound polypeptide. Hydrophobic amino acids include alanine, glycine, histidine, isoleucine, leucine, lysine, methionine, phenylalanine, proline, threonine, tryptophan, tyrosine, and valine

Identification of the Function of an Expression Product of a Full-Length Gene

Ribozymes, antisense constructs, and dominant negative mutants can be used to determine function of the expression product of a gene corresponding to a polynucleotide provided herein. These methods and compositions are particularly useful where the provided novel polynucleotide exhibits no significant or substantial homology to a sequence encoding a gene of known function. Antisense molecules and ribozymes can be constructed from synthetic polynucleotides. Typically, the phosphoramidite method of oligonucleotide synthesis is used. See Beaucage *et al.*, *Tet. Lett.* (1981) 22:1859 and USPN 4,668,777. Automated devices for synthesis are available to create oligonucleotides using this chemistry. Examples of such devices include Biosearch 8600, Models 392 and 394 by Applied Biosystems, a division of Perkin-Elmer Corp., Foster City, California, USA; and Expedite by Perceptive Biosystems, Framingham, Massachusetts, USA. Synthetic RNA, phosphate analog oligonucleotides, and chemically derivatized oligonucleotides can also be produced, and can be covalently attached to other molecules. RNA oligonucleotides can be synthesized, for example, using RNA phosphoramidites. This method can be performed on an automated synthesizer, such as Applied Biosystems, Models 392 and 394, Foster City, California, USA.

Phosphorothioate oligonucleotides can also be synthesized for antisense construction. A sulfurizing reagent, such as tetraethylthiuram disulfide (TETD) in acetonitrile can be used to convert the internucleotide cyanoethyl phosphite to the phosphorothioate triester within 15 minutes at room temperature. TETD replaces the iodine reagent, while all other reagents used for standard phosphoramidite chemistry remain the same. Such a synthesis method can be automated using Models 392 and 394 by Applied Biosystems, for example.

Oligonucleotides of up to 200 nt can be synthesized, more typically, 100 nt, more typically 50 nt; even more typically 30 to 40 nt. These synthetic fragments can be annealed and ligated together to construct larger fragments. See, for example, Sambrook *et al.*, *supra*. Trans-cleaving 10 catalytic RNAs (ribozymes) are RNA molecules possessing endoribonuclease activity. Ribozymes are specifically designed for a particular target, and the target message must contain a specific nucleotide sequence. They are engineered to cleave any RNA species site-specifically in the background of cellular RNA. The cleavage event renders the mRNA unstable and prevents protein expression. Importantly, ribozymes can be used to inhibit expression of a gene of unknown 15 function for the purpose of determining its function in an in vitro or in vivo context, by detecting the phenotypic effect. One commonly used ribozyme motif is the hammerhead, for which the substrate sequence requirements are minimal. Design of the hammerhead ribozyme, as well as therapeutic uses of ribozymes, are disclosed in Usman *et al.*, *Current Opin. Struct. Biol.* (1996) 6:527. Methods for production of ribozymes, including hairpin structure ribozyme fragments, 20 methods of increasing ribozyme specificity, and the like are known in the art.

The hybridizing region of the ribozyme can be modified or can be prepared as a branched structure as described in Horn and Urdea, *Nucleic Acids Res.* (1989) 17:6959. The basic structure of the ribozymes can also be chemically altered in ways familiar to those skilled in the art, and chemically synthesized ribozymes can be administered as synthetic oligonucleotide derivatives 25 modified by monomeric units. In a therapeutic context, liposome mediated delivery of ribozymes improves cellular uptake, as described in Birikh *et al.*, *Eur. J. Biochem.* (1997) 245:1.

Antisense nucleic acids are designed to specifically bind to RNA, resulting in the formation of RNA-DNA or RNA-RNA hybrids, with an arrest of DNA replication, reverse transcription or messenger RNA translation. Antisense polynucleotides based on a selected polynucleotide sequence 30 can interfere with expression of the corresponding gene. Antisense polynucleotides are typically generated within the cell by expression from antisense constructs that contain the antisense strand as the transcribed strand. Antisense polynucleotides based on the disclosed polynucleotides will bind and/or interfere with the translation of mRNA comprising a sequence complementary to the antisense polynucleotide. The expression products of control cells and cells treated with the 35 antisense construct are compared to detect the protein product of the gene corresponding to the

polynucleotide upon which the antisense construct is based. The protein is isolated and identified using routine biochemical methods.

Given the extensive background literature and clinical experience in antisense therapy, one skilled in the art can use selected polynucleotides of the invention as additional potential therapeutics. The choice of polynucleotide can be narrowed by first testing them for binding to "hot spot" regions of the genome of cancerous cells. If a polynucleotide is identified as binding to a "hot spot", testing the polynucleotide as an antisense compound in the corresponding cancer cells is warranted.

As an alternative method for identifying function of the gene corresponding to a polynucleotide disclosed herein, dominant negative mutations are readily generated for corresponding proteins that are active as homomultimers. A mutant polypeptide will interact with wild-type polypeptides (made from the other allele) and form a non-functional multimer. Thus, a mutation is in a substrate-binding domain, a catalytic domain, or a cellular localization domain. Preferably, the mutant polypeptide will be overproduced. Point mutations are made that have such an effect. In addition, fusion of different polypeptides of various lengths to the terminus of a protein can yield dominant negative mutants. General strategies are available for making dominant negative mutants (see, e.g., Herskowitz, *Nature* (1987) 329:219). Such techniques can be used to create loss of function mutations, which are useful for determining protein function.

Polypeptides and Variants Thereof

The polypeptides of the invention include those encoded by the disclosed polynucleotides, as well as nucleic acids that, by virtue of the degeneracy of the genetic code, are not identical in sequence to the disclosed polynucleotides. Thus, the invention includes within its scope a polypeptide encoded by a polynucleotide having the sequence of any one of SEQ ID NOS:1-2707 or a variant thereof.

In general, the term "polypeptide" as used herein refers to both the full length polypeptide encoded by the recited polynucleotide, the polypeptide encoded by the gene represented by the recited polynucleotide, as well as portions or fragments thereof. "Polypeptides" also includes variants of the naturally occurring proteins, where such variants are homologous or substantially similar to the naturally occurring protein, and can be of an origin of the same or different species as the naturally occurring protein (e.g., human, murine, or some other species that naturally expresses the recited polypeptide, usually a mammalian species). In general, variant polypeptides have a sequence that has at least about 80%, usually at least about 90%, and more usually at least about 98% sequence identity with a differentially expressed polypeptide of the invention, as measured by BLAST 2.0 using the parameters described above. The variant polypeptides can be naturally or non-

naturally glycosylated. *i.e.*, the polypeptide has a glycosylation pattern that differs from the glycosylation pattern found in the corresponding naturally occurring protein.

The invention also encompasses homologs of the disclosed polypeptides (or fragments thereof) where the homologs are isolated from other species. *i.e.*, other animal or plant species.

5 where such homologs, usually mammalian species, *e.g.*, rodents, such as mice, rats; domestic animals, *e.g.*, horse, cow, dog, cat; and humans. By "homolog" is meant a polypeptide having at least about 35%, usually at least about 40% and more usually at least about 60% amino acid sequence identity to a particular differentially expressed protein as identified above, where sequence identity is determined using the BLAST 2.0 algorithm, with the parameters described *supra*.

10 In general, the polypeptides of the subject invention are provided in a non-naturally occurring environment, *e.g.* are separated from their naturally occurring environment. In certain embodiments, the subject protein is present in a composition that is enriched for the protein as compared to a control. As such, purified polypeptide is provided, where by purified is meant that the protein is present in a composition that is substantially free of non-differentially expressed 15 polypeptides, where by substantially free is meant that less than 90%, usually less than 60% and more usually less than 50% of the composition is made up of non-differentially expressed polypeptides.

Also within the scope of the invention are variants: variants of polypeptides include mutants, fragments, and fusions. Mutants can include amino acid substitutions, additions or 20 deletions. The amino acid substitutions can be conservative amino acid substitutions or substitutions to eliminate non-essential amino acids, such as to alter a glycosylation site, a phosphorylation site or an acetylation site, or to minimize misfolding by substitution or deletion of one or more cysteine residues that are not necessary for function. Conservative amino acid substitutions are those that preserve the general charge, hydrophobicity/ hydrophilicity, and/or steric bulk of the amino acid 25 substituted. Variants can be designed so as to retain or have enhanced biological activity of a particular region of the protein (*e.g.*, a functional domain and/or, where the polypeptide is a member of a protein family, a region associated with a consensus sequence). Selection of amino acid alterations for production of variants can be based upon the accessibility (interior vs. exterior) of the amino acid (see, *e.g.*, Go *et al.* *Int. J. Peptide Protein Res.* (1980) 15:211), the thermostability of the 30 variant polypeptide (see, *e.g.*, Querol *et al.*, *Prot. Eng.* (1996) 9:265), desired glycosylation sites (see, *e.g.*, Olsen and Thomsen, *J. Gen. Microbiol.* (1991) 137:579), desired disulfide bridges (see, *e.g.*, Clarke *et al.*, *Biochemistry* (1993) 32:4322; and Wakarchuk *et al.*, *Protein Eng.* (1994) 7:1379), desired metal binding sites (see, *e.g.*, Toma *et al.*, *Biochemistry* (1991) 30:97, and Haezerbrouck *et al.*, *Protein Eng.* (1993) 6:643), and desired substitutions with in proline loops (see, *e.g.*, Masul *et*

al. Appl. Env. Microbiol. (1994) 60:3579). Cysteine-depleted muteins can be produced as disclosed in USPN 4,959,314.

Variants also include fragments of the polypeptides disclosed herein, particularly biologically active fragments and/or fragments corresponding to functional domains. Fragments of interest will typically be at least about 10 aa to at least about 15 aa in length, usually at least about 50 aa in length, and can be as long as 300 aa in length or longer, but will usually not exceed about 1000 aa in length, where the fragment will have a stretch of amino acids that is identical to a polypeptide encoded by a polynucleotide having a sequence of any SEQ ID NOS:1-2707, or a homolog thereof. The protein variants described herein are encoded by polynucleotides that are within the scope of the invention. The genetic code can be used to select the appropriate codons to construct the corresponding variants.

Computer-Related Embodiments

In general, a library of polynucleotides is a collection of sequence information, which information is provided in either biochemical form (e.g., as a collection of polynucleotide molecules), or in electronic form (e.g., as a collection of polynucleotide sequences stored in a computer-readable form, as in a computer system and/or as part of a computer program). The sequence information of the polynucleotides can be used in a variety of ways, e.g., as a resource for gene discovery, as a representation of sequences expressed in a selected cell type (e.g., cell type markers), and/or as markers of a given disease or disease state. In general, a disease marker is a representation of a gene product that is present in all cells affected by disease either at an increased or decreased level relative to a normal cell (e.g., a cell of the same or similar type that is not substantially affected by disease). For example, a polynucleotide sequence in a library can be a polynucleotide that represents an mRNA, polypeptide, or other gene product encoded by the polynucleotide, that is either overexpressed or underexpressed in a breast ductal cell affected by cancer relative to a normal (i.e., substantially disease-free) breast cell.

The nucleotide sequence information of the library can be embodied in any suitable form, e.g., electronic or biochemical forms. For example, a library of sequence information embodied in electronic form comprises an accessible computer data file (or, in biochemical form, a collection of nucleic acid molecules) that contains the representative nucleotide sequences of genes that are differentially expressed (e.g., overexpressed or underexpressed) as between, for example, i) a cancerous cell and a normal cell; ii) a cancerous cell and a dysplastic cell; iii) a cancerous cell and a cell affected by a disease or condition other than cancer; iv) a metastatic cancerous cell and a normal cell and/or non-metastatic cancerous cell; v) a malignant cancerous cell and a non-malignant cancerous cell (or a normal cell) and/or vi) a dysplastic cell relative to a normal cell. Other combinations and comparisons of cells affected by various diseases or stages of disease will be

readily apparent to the ordinarily skilled artisan. Biochemical embodiments of the library include a collection of nucleic acids that have the sequences of the genes in the library, where the nucleic acids can correspond to the entire gene in the library or to a fragment thereof, as described in greater detail below.

5 The polynucleotide libraries of the subject invention generally comprise sequence information of a plurality of polynucleotide sequences, where at least one of the polynucleotides has a sequence of any of SEQ ID NOS:1-2707. By plurality is meant at least 2, usually at least 3 and can include up to all of SEQ ID NOS:1-2707. The length and number of polynucleotides in the library will vary with the nature of the library, e.g., if the library is an oligonucleotide array, a cDNA 10 array, a computer database of the sequence information, etc.

Where the library is an electronic library, the nucleic acid sequence information can be present in a variety of media. "Media" refers to a manufacture, other than an isolated nucleic acid molecule, that contains the sequence information of the present invention. Such a manufacture provides the genome sequence or a subset thereof in a form that can be examined by means not directly applicable to the sequence as it exists in a nucleic acid. For example, the nucleotide sequence of the present invention, e.g. the nucleic acid sequences of any of the polynucleotides of SEQ ID NOS:1-2707, can be recorded on computer readable media, e.g. any medium that can be read and accessed directly by a computer. Such media include, but are not limited to: magnetic storage media, such as a floppy disc, a hard disc storage medium, and a magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories such as magnetic/optical storage media. One of skill in the art can readily appreciate how any of the presently known computer readable mediums can be used to create a manufacture comprising a recording of the present sequence information. "Recorded" refers to a process for storing information on computer readable medium, using any such methods as known in the art. Any convenient data storage structure can be chosen, based on the means used to access the stored information. A variety of data processor programs and formats can be used for storage, e.g. word processing text file, database format, etc. In addition to the sequence information, electronic versions of the libraries of the invention can be provided in conjunction or connection with other computer-readable information and/or other types of computer-readable files (e.g., searchable files, 20 executable files, etc, including, but not limited to, for example, search program software, etc.).

30 By providing the nucleotide sequence in computer readable form, the information can be accessed for a variety of purposes. Computer software to access sequence information is publicly available. For example, the gapped BLAST (Altschul *et al. Nucleic Acids Res.* (1997) 25:3389-3402) and BLAZE (Brutlag *et al. Comp. Chem.* (1993) 17:203) search algorithms on a Sybase

system can be used to identify open reading frames (ORFs) within the genome that contain homology to ORFs from other organisms.

As used herein, "a computer-based system" refers to the hardware means, software means, and data storage means used to analyze the nucleotide sequence information of the present invention. The minimum hardware of the computer-based systems of the present invention comprises a central processing unit (CPU), input means, output means, and data storage means. A skilled artisan can readily appreciate that any one of the currently available computer-based system are suitable for use in the present invention. The data storage means can comprise any manufacture comprising a recording of the present sequence information as described above, or a memory access means that can access such a manufacture.

"Search means" refers to one or more programs implemented on the computer-based system, to compare a target sequence or target structural motif, or expression levels of a polynucleotide in a sample, with the stored sequence information. Search means can be used to identify fragments or regions of the genome that match a particular target sequence or target motif. A variety of known algorithms are publicly known and commercially available, e.g. MacPattern (EMBL), BLASTN and BLASTX (NCBI). A "target sequence" can be any polynucleotide or amino acid sequence of six or more contiguous nucleotides or two or more amino acids, preferably from about 10 to 100 amino acids or from about 30 to 300 nt. A variety of comparing means can be used to accomplish comparison of sequence information from a sample (e.g., to analyze target sequences, target motifs, or relative expression levels) with the data storage means. A skilled artisan can readily recognize that any one of the publicly available homology search programs can be used as the search means for the computer based systems of the present invention to accomplish comparison of target sequences and motifs. Computer programs to analyze expression levels in a sample and in controls are also known in the art.

A "target structural motif," or "target motif," refers to any rationally selected sequence or combination of sequences in which the sequence(s) are chosen based on a three-dimensional configuration that is formed upon the folding of the target motif, or on consensus sequences of regulatory or active sites. There are a variety of target motifs known in the art. Protein target motifs include, but are not limited to, enzyme active sites and signal sequences. Nucleic acid target motifs include, but are not limited to, hairpin structures, promoter sequences and other expression elements such as binding sites for transcription factors.

A variety of structural formats for the input and output means can be used to input and output the information in the computer-based systems of the present invention. One format for an output means ranks the relative expression levels of different polynucleotides. Such presentation

provides a skilled artisan with a ranking of relative expression levels to determine a gene expression profile. .

As discussed above, the "library" of the invention also encompasses biochemical libraries of the polynucleotides of SEQ ID NOS:1-2707 . e.g., collections of nucleic acids representing the provided polynucleotides. The biochemical libraries can take a variety of forms, e.g., a solution of cDNAs, a pattern of probe nucleic acids stably associated with a surface of a solid support (i.e., an array) and the like. Of particular interest are nucleic acid arrays in which one or more of SEQ ID NOS:1-2707 is represented on the array. By array is meant a an article of manufacture that has at least a substrate with at least two distinct nucleic acid targets on one of its surfaces, where the number of distinct nucleic acids can be considerably higher, typically being at least 10 nt, usually at least 20 nt and often at least 25 nt. A variety of different array formats have been developed and are known to those of skill in the art. The arrays of the subject invention find use in a variety of applications, including gene expression analysis, drug screening, mutation analysis and the like, as disclosed in the above-listed exemplary patent documents.

In addition to the above nucleic acid libraries, analogous libraries of polypeptides are also provided, where the where the polypeptides of the library will represent at least a portion of the polypeptides encoded by SEQ ID NOS:1-2707.

Utilities

Use of Polynucleotide Probes in Mapping, and in Tissue Profiling

Polynucleotide probes, generally comprising at least 12 contiguous nt of a polynucleotide as shown in the Sequence Listing, are used for a variety of purposes, such as chromosome mapping of the polynucleotide and detection of transcription levels. Additional disclosure about preferred regions of the disclosed polynucleotide sequences is found in the Examples. A probe that hybridizes specifically to a polynucleotide disclosed herein should provide a detection signal at least 5-, 10-, or 20-fold higher than the background hybridization provided with other unrelated sequences.

Detection of Expression Levels. Nucleotide probes are used to detect expression of a gene corresponding to the provided polynucleotide. In Northern blots, mRNA is separated electrophoretically and contacted with a probe. A probe is detected as hybridizing to an mRNA species of a particular size. The amount of hybridization is quantitated to determine relative amounts of expression, for example under a particular condition. Probes are used for *in situ* hybridization to cells to detect expression. Probes can also be used *in vivo* for diagnostic detection of hybridizing sequences. Probes are typically labeled with a radioactive isotope. Other types of detectable labels can be used such as chromophores, fluors, and enzymes. Other examples of nucleotide hybridization assays are described in WO92/02526 and USPN 5,124,246.

Alternatively, the Polymerase Chain Reaction (PCR) is another means for detecting small amounts of target nucleic acids (see, e.g., Mullis *et al.*, *Meth. Enzymol.* (1987) 155:335; USPN 4,683,195; and USPN 4,683,202). Two primer polynucleotides nucleotides that hybridize with the target nucleic acids are used to prime the reaction. The primers can be composed of sequence within 5 or 3' and 5' to the polynucleotides of the Sequence Listing. Alternatively, if the primers are 3' and 5' to these polynucleotides, they need not hybridize to them or the complements. After amplification of the target with a thermostable polymerase, the amplified target nucleic acids can be detected by methods known in the art, e.g., Southern blot. mRNA or cDNA can also be detected by traditional blotting techniques (e.g., Southern blot, Northern blot, etc.) described in Sambrook *et al.*.

10 "Molecular Cloning: A Laboratory Manual" (New York, Cold Spring Harbor Laboratory, 1989) (e.g., without PCR amplification). In general, mRNA or cDNA generated from mRNA using a polymerase enzyme can be purified and separated using gel electrophoresis, and transferred to a solid support, such as nitrocellulose. The solid support is exposed to a labeled probe, washed to remove any unhybridized probe, and duplexes containing the labeled probe are detected.

15 Mapping. Polynucleotides of the present invention can be used to identify a chromosome on which the corresponding gene resides. Such mapping can be useful in identifying the function of the polynucleotide-related gene by its proximity to other genes with known function. Function can also be assigned to the polynucleotide-related gene when particular syndromes or diseases map to the same chromosome. For example, use of polynucleotide probes in identification and quantification

20 of nucleic acid sequence aberrations is described in USPN 5,783,387. An exemplary mapping method is fluorescence in situ hybridization (FISH), which facilitates comparative genomic hybridization to allow total genome assessment of changes in relative copy number of DNA sequences (see, e.g., Valdes *et al.*, *Methods in Molecular Biology* (1997) 68:1). Polynucleotides can also be mapped to particular chromosomes using, for example, radiation hybrids or

25 chromosome-specific hybrid panels. See Leach *et al.*, *Advances in Genetics*, (1995) 33:63-99; Walter *et al.*, *Nature Genetics* (1994) 7:22; Walter and Goodfellow, *Trends in Genetics* (1992) 9:352. Panels for radiation hybrid mapping are available from Research Genetics, Inc., Huntsville, Alabama, USA. Databases for markers using various panels are available via the world wide web at <http://F/shgc-www.stanford.edu>; and <http://www-genome.wi.mit.edu/cgi-bin/contig/rhMapper.pl>. The

30 statistical program RHMAP can be used to construct a map based on the data from radiation hybridization with a measure of the relative likelihood of one order versus another. RHMAP is available via the world wide web at <http://www.sph.umich.edu/group/statgen/software>. In addition, commercial programs are available for identifying regions of chromosomes commonly associated with disease, such as cancer.

Tissue Typing or Profiling. Expression of specific mRNA corresponding to the provided polynucleotides can vary in different cell types and can be tissue-specific. This variation of mRNA levels in different cell types can be exploited with nucleic acid probe assays to determine tissue types. For example, PCR, branched DNA probe assays, or blotting techniques utilizing nucleic acid probes substantially identical or complementary to polynucleotides listed in the Sequence Listing can determine the presence or absence of the corresponding cDNA or mRNA.

5 Tissue typing can be used to identify the developmental organ or tissue source of a metastatic lesion by identifying the expression of a particular marker of that organ or tissue. If a polynucleotide is expressed only in a specific tissue type, and a metastatic lesion is found to express 10 that polynucleotide, then the developmental source of the lesion has been identified. Expression of a particular polynucleotide can be assayed by detection of either the corresponding mRNA or the protein product. As would be readily apparent to any forensic scientist, the sequences disclosed herein are useful in differentiating human tissue from non-human tissue. In particular, these sequences are useful to differentiate human tissue from bird, reptile, and amphibian tissue, for 15 example.

20 Use of Polymorphisms. A polynucleotide of the invention can be used in forensics, genetic analysis, mapping, and diagnostic applications where the corresponding region of a gene is polymorphic in the human population. Any means for detecting a polymorphism in a gene can be used, including, but not limited to electrophoresis of protein polymorphic variants, differential sensitivity to restriction enzyme cleavage, and hybridization to allele-specific probes.

Antibody Production

25 Expression products of a polynucleotide of the invention, as well as the corresponding mRNA, cDNA, or complete gene, can be prepared and used for raising antibodies for experimental, diagnostic, and therapeutic purposes. For polynucleotides to which a corresponding gene has not been assigned, this provides an additional method of identifying the corresponding gene. The polynucleotide or related cDNA is expressed as described above, and antibodies are prepared. These antibodies are specific to an epitope on the polypeptide encoded by the polynucleotide, and can precipitate or bind to the corresponding native protein in a cell or tissue preparation or in a cell-free extract of an *in vitro* expression system.

30 Methods for production of antibodies that specifically bind a selected antigen are well known in the art. Immunogens for raising antibodies can be prepared by mixing a polypeptide encoded by a polynucleotide of the invention with an adjuvant, and/or by making fusion proteins with larger immunogenic proteins. Polypeptides can also be covalently linked to other larger immunogenic proteins, such as keyhole limpet hemocyanin. Immunogens are typically administered 35 intradermally, subcutaneously, or intramuscularly to experimental animals such as rabbits, sheep,

and mice, to generate antibodies. Monoclonal antibodies can be generated by isolating spleen cells and fusing myeloma cells to form hybridomas. Alternatively, the selected polynucleotide is administered directly, such as by intramuscular injection, and expressed in vivo. The expressed protein generates a variety of protein-specific immune responses, including 5 production of antibodies, comparable to administration of the protein.

Preparations of polyclonal and monoclonal antibodies specific for polypeptides encoded by a selected polynucleotide are made using standard methods known in the art. The antibodies specifically bind to epitopes present in the polypeptides encoded by polynucleotides disclosed in the Sequence Listing. Typically, at least 6, 8, 10, or 12 contiguous amino acids are required to form an 10 epitope. Epitopes that involve non-contiguous amino acids may require a longer polypeptide, e.g., at least 15, 25, or 50 amino acids. Antibodies that specifically bind to human polypeptides encoded by the provided polypeptides should provide a detection signal at least 5-, 10-, or 20-fold higher than a 15 detection signal provided with other proteins when used in Western blots or other immunochemical assays. Preferably, antibodies that specifically bind polypeptides of the invention do not bind to other proteins in immunochemical assays at detectable levels and can immunoprecipitate the specific polypeptide from solution.

The invention also contemplates naturally occurring antibodies specific for a polypeptide of the invention. For example, serum antibodies to a polypeptide of the invention in a human population can be purified by methods well known in the art, e.g., by passing antiserum over a 20 column to which the corresponding selected polypeptide or fusion protein is bound. The bound antibodies can then be eluted from the column, for example using a buffer with a high salt concentration.

In addition to the antibodies discussed above, the invention also contemplates genetically 25 engineered antibodies, antibody derivatives (e.g., single chain antibodies, antibody fragments (e.g., Fab, etc.)), according to methods well known in the art.

Polynucleotides or Arrays for Diagnostics

Polynucleotide arrays provide a high throughput technique that can assay a large number of 30 polynucleotide sequences in a sample. This technology can be used as a diagnostic and as a tool to test for differential expression, e.g., to determine function of an encoded protein. Arrays can be created by spotting polynucleotide probes onto a substrate (e.g., glass, nitrocellulose, etc.) in a two-dimensional matrix or array having bound probes. The probes can be bound to the substrate by either covalent bonds or by non-specific interactions, such as hydrophobic interactions. Samples of polynucleotides can be detectably labeled (e.g., using radioactive or fluorescent labels) and then hybridized to the probes. Double stranded polynucleotides, comprising the labeled sample 35 polynucleotides bound to probe polynucleotides, can be detected once the unbound portion of the

sample is washed away. Techniques for constructing arrays and methods of using these arrays are described in EP 799 897; WO 97/29212; WO 97/27317; EP 785 280; WO 97/02357; USPN 5,593,839; USPN 5,578,832; EP 728 520; USPN 5,599,695; EP 721 016; USPN 5,556,752; WO 95/22058; and USPN 5,631,734. Arrays can be used to, for example, examine differential expression of genes and can be used to determine gene function. For example, arrays can be used to detect differential expression of a polynucleotide between a test cell and control cell (e.g., cancer cells and normal cells). For example, high expression of a particular message in a cancer cell, which is not observed in a corresponding normal cell, can indicate a cancer specific gene product. Exemplary uses of arrays are further described in, for example, Pappalardo *et al.*, *Sem. Radiation Oncol.* (1998) 8:217; and Ramsay *Nature Biotechnol.* (1998) 16:40.

Differential Expression in Diagnosis

The polynucleotides of the invention can also be used to detect differences in expression levels between two cells, e.g., as a method to identify abnormal or diseased tissue in a human. For polynucleotides corresponding to profiles of protein families, the choice of tissue can be selected according to the putative biological function. In general, the expression of a gene corresponding to a specific polynucleotide is compared between a first tissue that is suspected of being diseased and a second, normal tissue of the human. The tissue suspected of being abnormal or diseased can be derived from a different tissue type of the human, but preferably it is derived from the same tissue type; for example an intestinal polyp or other abnormal growth should be compared with normal intestinal tissue. The normal tissue can be the same tissue as that of the test sample, or any normal tissue of the patient, especially those that express the polynucleotide-related gene of interest (e.g., brain, thymus, testis, heart, prostate, placenta, spleen, small intestine, skeletal muscle, pancreas, and the mucosal lining of the colon). A difference between the polynucleotide-related gene, mRNA, or protein in the two tissues which are compared, for example in molecular weight, amino acid or nucleotide sequence, or relative abundance, indicates a change in the gene, or a gene which regulates it, in the tissue of the human that was suspected of being diseased. Examples of detection of differential expression and its use in diagnosis of cancer are described in USPNs 5,688,641 and 5,677,125.

A genetic predisposition to disease in a human can also be detected by comparing expression levels of an mRNA or protein corresponding to a polynucleotide of the invention in a fetal tissue with levels associated in normal fetal tissue. Fetal tissues that are used for this purpose include, but are not limited to, amniotic fluid, chorionic villi, blood, and the blastomere of an in vitro-fertilized embryo. The comparable normal polynucleotide-related gene is obtained from any tissue. The mRNA or protein is obtained from a normal tissue of a human in which the polynucleotide-related gene is expressed. Differences such as alterations in the nucleotide sequence

or size of the same product of the fetal polynucleotide-related gene or mRNA, or alterations in the molecular weight, amino acid sequence, or relative abundance of fetal protein, can indicate a germline mutation in the polynucleotide-related gene of the fetus, which indicates a genetic predisposition to disease. In general, diagnostic, prognostic, and other methods of the invention based on differential expression involve detection of a level or amount of a gene product, particularly a differentially expressed gene product, in a test sample obtained from a patient suspected of having or being susceptible to a disease (e.g., breast cancer, lung cancer, colon cancer and/or metastatic forms thereof), and comparing the detected levels to those levels found in normal cells (e.g., cells substantially unaffected by cancer) and/or other control cells (e.g., to differentiate a cancerous cell from a cell affected by dysplasia). Furthermore, the severity of the disease can be assessed by comparing the detected levels of a differentially expressed gene product with those levels detected in samples representing the levels of differentially gene product associated with varying degrees of severity of disease. It should be noted that use of the term "diagnostic" herein is not necessarily meant to exclude "prognostic" or "prognosis," but rather is used as a matter of convenience.

The term "differentially expressed gene" is generally intended to encompass a polynucleotide that can, for example, include an open reading frame encoding a gene product (e.g., a polypeptide), and/or introns of such genes and adjacent 5' and 3' non-coding nucleotide sequences involved in the regulation of expression, up to about 20 kb beyond the coding region, but possibly further in either direction. The gene can be introduced into an appropriate vector for extrachromosomal maintenance or for integration into a host genome. In general, a difference in expression level associated with a decrease in expression level of at least about 25%, usually at least about 50% to 75%, more usually at least about 90% or more is indicative of a differentially expressed gene of interest, *i.e.*, a gene that is underexpressed or down-regulated in the test sample relative to a control sample. Furthermore, a difference in expression level associated with an increase in expression of at least about 25%, usually at least about 50% to 75%, more usually at least about 90% and can be at least about 1 ½-fold, usually at least about 2-fold to about 10-fold, and can be about 100-fold to about 1,000-fold increase relative to a control sample is indicative of a differentially expressed gene of interest, *i.e.*, an overexpressed or up-regulated gene.

"Differentially expressed polynucleotide" as used herein means a nucleic acid molecule (RNA or DNA) comprising a sequence that represents a differentially expressed gene. *e.g.*, the differentially expressed polynucleotide comprises a sequence (e.g., an open reading frame encoding a gene product) that uniquely identifies a differentially expressed gene so that detection of the differentially expressed polynucleotide in a sample is correlated with the presence of a differentially expressed gene in a sample. "Differentially expressed polynucleotides" is also meant to encompass

fragments of the disclosed polynucleotides, e.g., fragments retaining biological activity, as well as nucleic acids homologous, substantially similar, or substantially identical (e.g., having about 90% sequence identity) to the disclosed polynucleotides.

"Diagnosis" as used herein generally includes determination of a subject's susceptibility to a disease or disorder, determination as to whether a subject is presently affected by a disease or disorder, as well as to the prognosis of a subject affected by a disease or disorder (e.g., identification of pre-metastatic or metastatic cancerous states, stages of cancer, or responsiveness of cancer to therapy). The present invention particularly encompasses diagnosis of subjects in the context of breast cancer (e.g., carcinoma in situ (e.g., ductal carcinoma in situ), estrogen receptor (ER)-positive breast cancer, ER-negative breast cancer, or other forms and/or stages of breast cancer), lung cancer (e.g., small cell carcinoma, non-small cell carcinoma, mesothelioma, and other forms and/or stages of lung cancer), and colon cancer (e.g., adenomatous polyp, colorectal carcinoma, and other forms and/or stages of colon cancer).

"Sample" or "biological sample" as used throughout here are generally meant to refer to samples of biological fluids or tissues, particularly samples obtained from tissues, especially from cells of the type associated with the disease for which the diagnostic application is designed (e.g., ductal adenocarcinoma), and the like. "Samples" is also meant to encompass derivatives and fractions of such samples (e.g., cell lysates). Where the sample is solid tissue, the cells of the tissue can be dissociated or tissue sections can be analyzed.

Methods of the subject invention useful in diagnosis or prognosis typically involve comparison of the abundance of a selected differentially expressed gene product in a sample of interest with that of a control to determine any relative differences in the expression of the gene product, where the difference can be measured qualitatively and/or quantitatively. Quantitation can be accomplished, for example, by comparing the level of expression product detected in the sample with the amounts of product present in a standard curve. A comparison can be made visually: by using a technique such as densitometry, with or without computerized assistance; by preparing a representative library of cDNA clones of mRNA isolated from a test sample, sequencing the clones in the library to determine that number of cDNA clones corresponding to the same gene product, and analyzing the number of clones corresponding to that same gene product relative to the number of clones of the same gene product in a control sample; or by using an array to detect relative levels of hybridization to a selected sequence or set of sequences, and comparing the hybridization pattern to that of a control. The differences in expression are then correlated with the presence or absence of an abnormal expression pattern. A variety of different methods for determining the nucleic acid abundance in a sample are known to those of skill in the art (see, e.g., WO 97/27317). In general, diagnostic assays of the invention involve detection of a gene product of a the polynucleotide

sequence (e.g., mRNA or polypeptide) that corresponds to a sequence of SEQ ID NOS:1-2707. The patient from whom the sample is obtained can be apparently healthy, susceptible to disease (e.g., as determined by family history or exposure to certain environmental factors), or can already be identified as having a condition in which altered expression of a gene product of the invention is 5 implicated.

Diagnosis can be determined based on detected gene product expression levels of a gene product encoded by at least one, preferably at least two or more, at least 3 or more, or at least 4 or more of the polynucleotides having a sequence set forth in SEQ ID NOS:1-2707, and can involve detection of expression of genes corresponding to all of SEQ ID NOS:1-2707 and/or additional 10 sequences that can serve as additional diagnostic markers and/or reference sequences. Where the diagnostic method is designed to detect the presence or susceptibility of a patient to cancer, the assay preferably involves detection of a gene product encoded by a gene corresponding to a polynucleotide that is differentially expressed in cancer. Examples of such differentially expressed polynucleotides are described in the Examples below. Given the provided polynucleotides and 15 information regarding their relative expression levels provided herein, assays using such polynucleotides and detection of their expression levels in diagnosis and prognosis will be readily apparent to the ordinarily skilled artisan.

Any of a variety of detectable labels can be used in connection with the various 20 embodiments of the diagnostic methods of the invention. Suitable detectable labels include fluorochromes, (e.g. fluorescein isothiocyanate (FITC), rhodamine, Texas Red, phycoerythrin, allophycocyanin, 6-carboxyfluorescein (6-FAM), 2',7'-dimethoxy-4',5'-dichloro-6-carboxyfluorescein, 6-carboxy-X-rhodamine (ROX), 6-carboxy-2',4',7',4,7-hexachlorofluorescein (HEX), 5-carboxyfluorescein (5-FAM) or N,N,N',N'-tetramethyl-6-carboxyrhodamine (TAMRA)), radioactive labels, (e.g. ^{32}P , ^{35}S , ^{3}H , etc.), and the like. The detectable label can involve a two 25 stage systems (e.g., biotin-avidin, hapten-anti-hapten antibody, etc.)

Reagents specific for the polynucleotides and polypeptides of the invention, such as 30 antibodies and nucleotide probes, can be supplied in a kit for detecting the presence of an expression product in a biological sample. The kit can also contain buffers or labeling components, as well as instructions for using the reagents to detect and quantify expression products in the biological sample. Exemplary embodiments of the diagnostic methods of the invention are described below in more detail.

Polypeptide detection in diagnosis. In one embodiment, the test sample is assayed for the 35 level of a differentially expressed polypeptide. Diagnosis can be accomplished using any of a number of methods to determine the absence or presence or altered amounts of the differentially expressed polypeptide in the test sample. For example, detection can utilize staining of cells or

histological sections with labeled antibodies, performed in accordance with conventional methods. Cells can be permeabilized to stain cytoplasmic molecules. In general, antibodies that specifically bind a differentially expressed polypeptide of the invention are added to a sample, and incubated for a period of time sufficient to allow binding to the epitope, usually at least about 10 minutes. The 5 antibody can be detectably labeled for direct detection (e.g., using radioisotopes, enzymes, fluorescers, chemiluminescers, and the like), or can be used in conjunction with a second stage antibody or reagent to detect binding (e.g., biotin with horseradish peroxidase-conjugated avidin, a secondary antibody conjugated to a fluorescent compound, e.g. fluorescein, rhodamine, Texas red, etc.). The absence or presence of antibody binding can be determined by various methods, including 10 flow cytometry of dissociated cells, microscopy, radiography, scintillation counting, etc. Any suitable alternative methods can of qualitative or quantitative detection of levels or amounts of differentially expressed polypeptide can be used, for example ELISA, western blot, immunoprecipitation, radioimmunoassay, etc.

mRNA detection. The diagnostic methods of the invention can also or alternatively involve 15 detection of mRNA encoded by a gene corresponding to a differentially expressed polynucleotides of the invention. Any suitable qualitative or quantitative methods known in the art for detecting specific mRNAs can be used. mRNA can be detected by, for example, *in situ* hybridization in tissue sections, by reverse transcriptase-PCR, or in Northern blots containing poly A+ mRNA. One of skill in the art can readily use these methods to determine differences in the size or amount of mRNA 20 transcripts between two samples. mRNA expression levels in a sample can also be determined by generation of a library of expressed sequence tags (ESTs) from the sample, where the EST library is representative of sequences present in the sample (Adams, et al., (1991) *Science* 252:1651). Enumeration of the relative representation of ESTs within the library can be used to approximate the 25 relative representation of the gene transcript within the starting sample. The results of EST analysis of a test sample can then be compared to EST analysis of a reference sample to determine the relative expression levels of a selected polynucleotide, particularly a polynucleotide corresponding to one or more of the differentially expressed genes described herein. Alternatively, gene expression in a test sample can be performed using serial analysis of gene expression (SAGE) methodology (e.g., Velculescu et al., *Science* (1995) 270:484) or differential display (DD) methodology (see, e.g., 30 U.S. 5,776,683; and U.S. 5,807,680).

Alternatively, gene expression can be analyzed using hybridization analysis. Oligonucleotides or cDNA can be used to selectively identify or capture DNA or RNA of specific sequence composition, and the amount of RNA or cDNA hybridized to a known capture sequence determined qualitatively or quantitatively, to provide information about the relative representation of 35 a particular message within the pool of cellular messages in a sample. Hybridization analysis can be

designed to allow for concurrent screening of the relative expression of hundreds to thousands of genes by using, for example, array-based technologies having high density formats, including filters, microscope slides, or microchips, or solution-based technologies that use spectroscopic analysis (e.g., mass spectrometry). One exemplary use of arrays in the diagnostic methods of the invention is 5 described below in more detail.

Use of a single gene in diagnostic applications. The diagnostic methods of the invention can focus on the expression of a single differentially expressed gene. For example, the diagnostic method can involve detecting a differentially expressed gene, or a polymorphism of such a gene (e.g., a polymorphism in a coding region or control region), that is associated with disease.

10 Disease-associated polymorphisms can include deletion or truncation of the gene, mutations that alter expression level and/or affect activity of the encoded protein, *etc.*

A number of methods are available for analyzing nucleic acids for the presence of a specific sequence, *e.g.* a disease associated polymorphism. Where large amounts of DNA are available, genomic DNA is used directly. Alternatively, the region of interest is cloned into a suitable vector 15 and grown in sufficient quantity for analysis. Cells that express a differentially expressed gene can be used as a source of mRNA, which can be assayed directly or reverse transcribed into cDNA for analysis. The nucleic acid can be amplified by conventional techniques, such as the polymerase chain reaction (PCR), to provide sufficient amounts for analysis, and a detectable label can be included in the amplification reaction (*e.g.*, using a detectably labeled primer or detectably labeled 20 oligonucleotides) to facilitate detection. Alternatively, various methods are also known in the art that utilize oligonucleotide ligation as a means of detecting polymorphisms, *see e.g.*, Riley *et al.*, *Nucl. Acids Res.* (1990) 18:2887; and Delahunt *et al.*, *Am. J. Hum. Genet.* (1996) 58:1239.

The amplified or cloned sample nucleic acid can be analyzed by one of a number of methods known in the art. The nucleic acid can be sequenced by dideoxy or other methods, and the sequence 25 of bases compared to a selected sequence, *e.g.*, to a wild-type sequence. Hybridization with the polymorphic or variant sequence can also be used to determine its presence in a sample (*e.g.*, by Southern blot, dot blot, *etc.*). The hybridization pattern of a polymorphic or variant sequence and a control sequence to an array of oligonucleotide probes immobilized on a solid support, as described 30 in US 5,445,934, or in WO 95/35505, can also be used as a means of identifying polymorphic or variant sequences associated with disease. Single strand conformational polymorphism (SSCP) analysis, denaturing gradient gel electrophoresis (DGGE), and heteroduplex analysis in gel matrices are used to detect conformational changes created by DNA sequence variation as alterations in electrophoretic mobility. Alternatively, where a polymorphism creates or destroys a recognition site 35 for a restriction endonuclease, the sample is digested with that endonuclease, and the products size

fractionated to determine whether the fragment was digested. Fractionation is performed by gel or capillary electrophoresis, particularly acrylamide or agarose gels.

Screening for mutations in a gene can be based on the functional or antigenic characteristics of the protein. Protein truncation assays are useful in detecting deletions that can affect the 5 biological activity of the protein. Various immunoassays designed to detect polymorphisms in proteins can be used in screening. Where many diverse genetic mutations lead to a particular disease phenotype, functional protein assays have proven to be effective screening tools. The activity of the encoded protein can be determined by comparison with the wild-type protein.

Pattern matching in diagnosis using arrays. In another embodiment, the diagnostic and/or 10 prognostic methods of the invention involve detection of expression of a selected set of genes in a test sample to produce a test expression pattern (TEP). The TEP is compared to a reference expression pattern (REP), which is generated by detection of expression of the selected set of genes in a reference sample (e.g., a positive or negative control sample). The selected set of genes includes at least one of the genes of the invention, which genes correspond to the polynucleotide 15 sequences of SEQ ID NOS:1-2707. Of particular interest is a selected set of genes that includes gene differentially expressed in the disease for which the test sample is to be screened.

"Reference sequences" or "reference polynucleotides" as used herein in the context of differential gene expression analysis and diagnosis/prognosis refers to a selected set of polynucleotides, which selected set includes at least one or more of the differentially expressed 20 polynucleotides described herein. A plurality of reference sequences, preferably comprising positive and negative control sequences, can be included as reference sequences. Additional suitable reference sequences are found in GenBank, Unigene, and other nucleotide sequence databases (including, e.g., expressed sequence tag (EST), partial, and full-length sequences).

"Reference array" means an array having reference sequences for use in hybridization with a 25 sample, where the reference sequences include all, at least one of, or any subset of the differentially expressed polynucleotides described herein. Usually such an array will include at least 3 different reference sequences, and can include any one or all of the provided differentially expressed sequences. Arrays of interest can further comprise sequences, including polymorphisms, of other 30 genetic sequences, particularly other sequences of interest for screening for a disease or disorder (e.g., cancer, dysplasia, or other related or unrelated diseases, disorders, or conditions). The oligonucleotide sequence on the array will usually be at least about 12 nt in length, and can be of about the length of the provided sequences, or can extend into the flanking regions to generate 35 fragments of 100 nt to 200 nt in length or more. Reference arrays can be produced according to any suitable methods known in the art. For example, methods of producing large arrays of oligonucleotides are described in U.S. 5,134,854, and U.S. 5,445,934 using light-directed synthesis

techniques. Using a computer controlled system, a heterogeneous array of monomers is converted, through simultaneous coupling at a number of reaction sites, into a heterogeneous array of polymers. Alternatively, microarrays are generated by deposition of pre-synthesized oligonucleotides onto a solid substrate, for example as described in PCT published application no. WO 95/35505.

5 A "reference expression pattern" or "REP" as used herein refers to the relative levels of expression of a selected set of genes, particularly of differentially expressed genes, that is associated with a selected cell type, *e.g.*, a normal cell, a cancerous cell, a cell exposed to an environmental stimulus, and the like. A "test expression pattern" or "TEP" refers to relative levels of expression of a selected set of genes, particularly of differentially expressed genes, in a test sample (*e.g.*, a cell of 10 unknown or suspected disease state, from which mRNA is isolated).

REPs can be generated in a variety of ways according to methods well known in the art. For example, REPs can be generated by hybridizing a control sample to an array having a selected set of polynucleotides (particularly a selected set of differentially expressed polynucleotides), acquiring the hybridization data from the array, and storing the data in a format that allows for ready 15 comparison of the REP with a TEP. Alternatively, all expressed sequences in a control sample can be isolated and sequenced, *e.g.*, by isolating mRNA from a control sample, converting the mRNA into cDNA, and sequencing the cDNA. The resulting sequence information roughly or precisely reflects the identity and relative number of expressed sequences in the sample. The sequence information can then be stored in a format (*e.g.*, a computer-readable format) that allows for ready 20 comparison of the REP with a TEP. The REP can be normalized prior to or after data storage, and/or can be processed to selectively remove sequences of expressed genes that are of less interest or that might complicate analysis (*e.g.*, some or all of the sequences associated with housekeeping genes can be eliminated from REP data).

TEPs can be generated in a manner similar to REPs, *e.g.*, by hybridizing a test sample to an 25 array having a selected set of polynucleotides, particularly a selected set of differentially expressed polynucleotides, acquiring the hybridization data from the array, and storing the data in a format that allows for ready comparison of the TEP with a REP. The REP and TEP to be used in a comparison can be generated simultaneously, or the TEP can be compared to previously generated and stored REPs.

30 In one embodiment of the invention, comparison of a TEP with a REP involves hybridizing a test sample with a reference array, where the reference array has one or more reference sequences for use in hybridization with a sample. The reference sequences include all, at least one of, or any subset of the differentially expressed polynucleotides described herein. Hybridization data for the test sample is acquired, the data normalized, and the produced TEP compared with a REP generated 35 using an array having the same or similar selected set of differentially expressed polynucleotides.

Probes that correspond to sequences differentially expressed between the two samples will show decreased or increased hybridization efficiency for one of the samples relative to the other.

Methods for collection of data from hybridization of samples with a reference arrays are well known in the art. For example, the polynucleotides of the reference and test samples can be

5 generated using a detectable fluorescent label, and hybridization of the polynucleotides in the samples detected by scanning the microarrays for the presence of the detectable label using, for example, a microscope and light source for directing light at a substrate. A photon counter detects fluorescence from the substrate, while an x-y translation stage varies the location of the substrate. A confocal detection device that can be used in the subject methods is described in USPN 5,631,734.

10 A scanning laser microscope is described in Shalon et al., *Genome Res.* (1996) 6:639. A scan, using the appropriate excitation line, is performed for each fluorophore used. The digital images generated from the scan are then combined for subsequent analysis. For any particular array element, the ratio of the fluorescent signal from one sample (e.g., a test sample) is compared to the fluorescent signal from another sample (e.g., a reference sample), and the relative signal intensity

15 determined.

Methods for analyzing the data collected from hybridization to arrays are well known in the art. For example, where detection of hybridization involves a fluorescent label, data analysis can include the steps of determining fluorescent intensity as a function of substrate position from the data collected, removing outliers, i.e. data deviating from a predetermined statistical distribution,

20 and calculating the relative binding affinity of the targets from the remaining data. The resulting data can be displayed as an image with the intensity in each region varying according to the binding affinity between targets and probes.

In general, the test sample is classified as having a gene expression profile corresponding to that associated with a disease or non-disease state by comparing the TEP generated from the test

25 sample to one or more REPs generated from reference samples (e.g., from samples associated with cancer or specific stages of cancer, dysplasia, samples affected by a disease other than cancer, normal samples, etc.). The criteria for a match or a substantial match between a TEP and a REP include expression of the same or substantially the same set of reference genes, as well as expression of these reference genes at substantially the same levels (e.g., no significant difference between the

30 samples for a signal associated with a selected reference sequence after normalization of the samples, or at least no greater than about 25% to about 40% difference in signal strength for a given reference sequence. In general, a pattern match between a TEP and a REP includes a match in expression, preferably a match in qualitative or quantitative expression level, of at least one of, all or any subset of the differentially expressed genes of the invention.

Pattern matching can be performed manually, or can be performed using a computer program. Methods for preparation of substrate matrices (e.g., arrays), design of oligonucleotides for use with such matrices, labeling of probes, hybridization conditions, scanning of hybridized matrices, and analysis of patterns generated, including comparison analysis, are described in, for 5 example, U.S. 5,800,992.

Diagnosis, Prognosis and Management of Cancer

The polynucleotides of the invention and their gene products are of particular interest as genetic or biochemical markers (e.g., in blood or tissues) that will detect the earliest changes along the carcinogenesis pathway and/or to monitor the efficacy of various therapies and preventive 10 interventions. For example, the level of expression of certain polynucleotides can be indicative of a poorer prognosis, and therefore warrant more aggressive chemo- or radio-therapy for a patient or vice versa. The correlation of novel surrogate tumor specific features with response to treatment and outcome in patients can define prognostic indicators that allow the design of tailored therapy based on the molecular profile of the tumor. These therapies include antibody targeting and gene therapy. 15 Determining expression of certain polynucleotides and comparison of a patients profile with known expression in normal tissue and variants of the disease allows a determination of the best possible treatment for a patient, both in terms of specificity of treatment and in terms of comfort level of the patient. Surrogate tumor markers, such as polynucleotide expression, can also be used to better classify, and thus diagnose and treat, different forms and disease states of cancer. Two 20 classifications widely used in oncology that can benefit from identification of the expression levels of the polynucleotides of the invention are staging of the cancerous disorder, and grading the nature of the cancerous tissue.

The polynucleotides of the invention can be useful to monitor patients having or susceptible to cancer to detect potentially malignant events at a molecular level before they are detectable at a 25 gross morphological level. Furthermore, a polynucleotide of the invention identified as important for one type of cancer can also have implications for development or risk of development of other types of cancer, e.g., where a polynucleotide is differentially expressed across various cancer types. Thus, for example, expression of a polynucleotide that has clinical implications for metastatic colon cancer can also have clinical implications for stomach cancer or endometrial cancer.

30 Staging. Staging is a process used by physicians to describe how advanced the cancerous state is in a patient. Staging assists the physician in determining a prognosis, planning treatment and evaluating the results of such treatment. Staging systems vary with the types of cancer, but generally involve the following "TNM" system: the type of tumor, indicated by T; whether the cancer has metastasized to nearby lymph nodes, indicated by N; and whether the cancer has metastasized to 35 more distant parts of the body, indicated by M. Generally, if a cancer is only detectable in the area

of the primary lesion without having spread to any lymph nodes it is called Stage I. If it has spread only to the closest lymph nodes, it is called Stage II. In Stage III, the cancer has generally spread to the lymph nodes in near proximity to the site of the primary lesion. Cancers that have spread to a distant part of the body, such as the liver, bone, brain or other site, are Stage IV, the most advanced 5 stage.

The polynucleotides of the invention can facilitate fine-tuning of the staging process by identifying markers for the aggressivity of a cancer, e.g. the metastatic potential, as well as the presence in different areas of the body. Thus, a Stage II cancer with a polynucleotide signifying a high metastatic potential cancer can be used to change a borderline Stage II tumor to a Stage III 10 tumor, justifying more aggressive therapy. Conversely, the presence of a polynucleotide signifying a lower metastatic potential allows more conservative staging of a tumor.

Grading of cancers. Grade is a term used to describe how closely a tumor resembles normal tissue of its same type. The microscopic appearance of a tumor is used to identify tumor grade based on parameters such as cell morphology, cellular organization, and other markers of differentiation.

15 As a general rule, the grade of a tumor corresponds to its rate of growth or aggressiveness, with undifferentiated or high-grade tumors being more aggressive than well differentiated or low-grade tumors. The following guidelines are generally used for grading tumors: 1) GX Grade cannot be assessed; 2) G1 Well differentiated; G2 Moderately well differentiated; 3) G3 Poorly differentiated; 4) G4 Undifferentiated. The polynucleotides of the invention can be especially valuable in 20 determining the grade of the tumor, as they not only can aid in determining the differentiation status of the cells of a tumor, they can also identify factors other than differentiation that are valuable in determining the aggressiveness of a tumor, such as metastatic potential.

Detection of lung cancer. The polynucleotides of the invention can be used to detect lung 25 cancer in a subject. Although there are more than a dozen different kinds of lung cancer, the two main types of lung cancer are small cell and nonsmall cell, which encompass about 90% of all lung cancer cases. Small cell carcinoma (also called oat cell carcinoma) usually starts in one of the larger bronchial tubes, grows fairly rapidly, and is likely to be large by the time of diagnosis. Nonsmall cell lung cancer (NSCLC) is made up of three general subtypes of lung cancer. Epidermoid carcinoma (also called squamous cell carcinoma) usually starts in one of the larger bronchial tubes 30 and grows relatively slowly. The size of these tumors can range from very small to quite large. Adenocarcinoma starts growing near the outside surface of the lung and can vary in both size and growth rate. Some slowly growing adenocarcinomas are described as alveolar cell cancer. Large cell carcinoma starts near the surface of the lung, grows rapidly, and the growth is usually fairly large when diagnosed. Other less common forms of lung cancer are carcinoid, cylindroma, 35 mucoepidermoid, and malignant mesothelioma.

The polynucleotides of the invention, e.g., polynucleotides differentially expressed in normal cells versus cancerous lung cells (e.g., tumor cells of high or low metastatic potential) or between types of cancerous lung cells (e.g., high metastatic versus low metastatic), can be used to distinguish types of lung cancer as well as identifying traits specific to a certain patient's cancer and 5 selecting an appropriate therapy. For example, if the patient's biopsy expresses a polynucleotide that is associated with a low metastatic potential, it may justify leaving a larger portion of the patient's lung in surgery to remove the lesion. Alternatively, a smaller lesion with expression of a polynucleotide that is associated with high metastatic potential may justify a more radical removal of lung tissue and/or the surrounding lymph nodes, even if no metastasis can be identified through 10 pathological examination.

Detection of breast cancer. The majority of breast cancers are adenocarcinomas subtypes, which can be summarized as follows: 1) ductal carcinoma in situ (DCIS), including comedocarcinoma; 2) infiltrating (or invasive) ductal carcinoma (IDC); 3) lobular carcinoma in situ (LCIS); 4) infiltrating (or invasive) lobular carcinoma (ILC); 5) inflammatory breast cancer; 6) 15 medullary carcinoma; 7) mucinous carcinoma; 8) Paget's disease of the nipple; 9) Phyllodes tumor; and 10) tubular carcinoma;

The expression of polynucleotides of the invention can be used in the diagnosis and management of breast cancer, as well as to distinguish between types of breast cancer. Detection of breast cancer can be determined using expression levels of any of the appropriate polynucleotides of 20 the invention, either alone or in combination. Determination of the aggressive nature and/or the metastatic potential of a breast cancer can also be determined by comparing levels of one or more polynucleotides of the invention and comparing levels of another sequence known to vary in cancerous tissue, e.g. ER expression. In addition, development of breast cancer can be detected by examining the ratio of expression of a differentially expressed polynucleotide to the levels of steroid 25 hormones (e.g., testosterone or estrogen) or to other hormones (e.g., growth hormone, insulin). Thus expression of specific marker polynucleotides can be used to discriminate between normal and cancerous breast tissue, to discriminate between breast cancers with different cells of origin, to discriminate between breast cancers with different potential metastatic rates, etc.

Detection of colon cancer. The polynucleotides of the invention exhibiting the appropriate 30 expression pattern can be used to detect colon cancer in a subject. Colorectal cancer is one of the most common neoplasms in humans and perhaps the most frequent form of hereditary neoplasia. Prevention and early detection are key factors in controlling and curing colorectal cancer. Colorectal cancer begins as polyps, which are small, benign growths of cells that form on the inner 35 lining of the colon. Over a period of several years, some of these polyps accumulate additional mutations and become cancerous. Multiple familial colorectal cancer disorders have been identified.

which are summarized as follows: 1) Familial adenomatous polyposis (FAP); 2) Gardner's syndrome; 3) Hereditary nonpolyposis colon cancer (HNPCC); and 4) Familial colorectal cancer in Ashkenazi Jews. The expression of appropriate polynucleotides of the invention can be used in the diagnosis, prognosis and management of colorectal cancer. Detection of colon cancer can be
5 determined using expression levels of any of these sequences alone or in combination with the levels of expression. Determination of the aggressive nature and/or the metastatic potential of a colon cancer can be determined by comparing levels of one or more polynucleotides of the invention and comparing total levels of another sequence known to vary in cancerous tissue, e.g., expression of p53, DCC ras, for FAP (see, e.g., Fearon ER, *et al.*, *Cell* (1990) 61(5):759; Hamilton SR *et al.*,
10 *Cancer* (1993) 72:957; Bodmer W, *et al.*, *Nat Genet.* (1994) 4(3):217; Fearon ER, *Ann N Y Acad Sci.* (1995) 768:101). For example, development of colon cancer can be detected by examining the ratio of any of the polynucleotides of the invention to the levels of oncogenes (e.g. ras) or tumor suppressor genes (e.g. FAP or p53). Thus expression of specific marker polynucleotides can be used to discriminate between normal and cancerous colon tissue, to discriminate between colon cancers
15 with different cells of origin, to discriminate between colon cancers with different potential metastatic rates, etc.

Use of Polynucleotides to Screen for Peptide Analogs and Antagonists

Polypeptides encoded by the instant polynucleotides and corresponding full length genes can be used to screen peptide libraries to identify binding partners, such as receptors, from among
20 the encoded polypeptides. Peptide libraries can be synthesized according to methods known in the art (see, e.g., USPN 5,010,175, and WO 91/17823). Agonists or antagonists of the polypeptides of the invention can be screened using any available method known in the art, such as signal transduction, antibody binding, receptor binding, mitogenic assays, chemotaxis assays, etc. The assay conditions ideally should resemble the conditions under which the native activity is exhibited
25 *in vivo*, that is, under physiologic pH, temperature, and ionic strength. Suitable agonists or antagonists will exhibit strong inhibition or enhancement of the native activity at concentrations that do not cause toxic side effects in the subject. Agonists or antagonists that compete for binding to the native polypeptide can require concentrations equal to or greater than the native concentration, while inhibitors capable of binding irreversibly to the polypeptide can be added in concentrations on the
30 order of the native concentration.

Such screening and experimentation can lead to identification of a novel polypeptide binding partner, such as a receptor, encoded by a gene or a cDNA corresponding to a polynucleotide of the invention, and at least one peptide agonist or antagonist of the novel binding partner. Such agonists and antagonists can be used to modulate, enhance, or inhibit receptor function in cells to
35 which the receptor is native, or in cells that possess the receptor as a result of genetic engineering.

Further, if the novel receptor shares biologically important characteristics with a known receptor, information about agonist/antagonist binding can facilitate development of improved agonists/antagonists of the known receptor.

Pharmaceutical Compositions and Therapeutic Uses

5 Pharmaceutical compositions of the invention can comprise polypeptides, antibodies, or polynucleotides (including antisense nucleotides and ribozymes) of the claimed invention in a therapeutically effective amount. The term "therapeutically effective amount" as used herein refers to an amount of a therapeutic agent to treat, ameliorate, or prevent a desired disease or condition, or to exhibit a detectable therapeutic or preventative effect. The effect can be detected by, for example, 10 chemical markers or antigen levels. Therapeutic effects also include reduction in physical symptoms, such as decreased body temperature. The precise effective amount for a subject will depend upon the subject's size and health, the nature and extent of the condition, and the therapeutics or combination of therapeutics selected for administration. Thus, it is not useful to specify an exact effective amount in advance. However, the effective amount for a given situation is determined by 15 routine experimentation and is within the judgment of the clinician. For purposes of the present invention, an effective dose will generally be from about 0.01 mg/kg to 50 mg/kg or 0.05 mg/kg to about 10 mg/kg of the DNA constructs in the individual to which it is administered.

A pharmaceutical composition can also contain a pharmaceutically acceptable carrier. The term "pharmaceutically acceptable carrier" refers to a carrier for administration of a therapeutic 20 agent, such as antibodies or a polypeptide, genes, and other therapeutic agents. The term refers to any pharmaceutical carrier that does not itself induce the production of antibodies harmful to the individual receiving the composition, and which can be administered without undue toxicity.

Suitable carriers can be large, slowly metabolized macromolecules such as proteins, 25 polysaccharides, polylactic acids, polyglycolic acids, polymeric amino acids, amino acid copolymers, and inactive virus particles. Such carriers are well known to those of ordinary skill in the art. Pharmaceutically acceptable carriers in therapeutic compositions can include liquids such as water, saline, glycerol and ethanol. Auxiliary substances, such as wetting or emulsifying agents, pH buffering substances, and the like, can also be present in such vehicles. Typically, the therapeutic compositions are prepared as injectables, either as liquid solutions or suspensions; solid forms 30 suitable for solution in, or suspension in, liquid vehicles prior to injection can also be prepared. Liposomes are included within the definition of a pharmaceutically acceptable carrier. Pharmaceutically acceptable salts can also be present in the pharmaceutical composition, e.g., mineral acid salts such as hydrochlorides, hydrobromides, phosphates, sulfates, and the like; and the salts of organic acids such as acetates, propionates, malonates, benzoates, and the like. A thorough

discussion of pharmaceutically acceptable excipients is available in *Remington's Pharmaceutical Sciences* (Mack Pub. Co., N.J. 1991).

Delivery Methods. Once formulated, the compositions of the invention can be (1) administered directly to the subject (e.g., as polynucleotide or polypeptides); or (2) delivered ex vivo, to cells derived from the subject (e.g., as in *ex vivo* gene therapy). Direct delivery of the compositions will generally be accomplished by parenteral injection, e.g., subcutaneously, intraperitoneally, intravenously or intramuscularly, intratumoral or to the interstitial space of a tissue. Other modes of administration include oral and pulmonary administration, suppositories, and transdermal applications, needles, and gene guns or hyposprays. Dosage treatment can be a single dose schedule or a multiple dose schedule.

Methods for the *ex vivo* delivery and reimplantation of transformed cells into a subject are known in the art and described in e.g., International Publication No. WO 93/14778. Examples of cells useful in *ex vivo* applications include, for example, stem cells, particularly hematopoietic, lymph cells, macrophages, dendritic cells, or tumor cells. Generally, delivery of nucleic acids for both *ex vivo* and *in vitro* applications can be accomplished by, for example, dextran-mediated transfection, calcium phosphate precipitation, polybrene mediated transfection, protoplast fusion, electroporation, encapsulation of the polynucleotide(s) in liposomes, and direct microinjection of the DNA into nuclei, all well known in the art.

Once a gene corresponding to a polynucleotide of the invention has been found to correlate with a proliferative disorder, such as neoplasia, dysplasia, and hyperplasia, the disorder can be amenable to treatment by administration of a therapeutic agent based on the provided polynucleotide, corresponding polypeptide or other corresponding molecule (e.g., antisense, ribozyme, etc.).

The dose and the means of administration of the inventive pharmaceutical compositions are determined based on the specific qualities of the therapeutic composition, the condition, age, and weight of the patient, the progression of the disease, and other relevant factors. For example, administration of polynucleotide therapeutic compositions agents of the invention includes local or systemic administration, including injection, oral administration, particle gun or catheterized administration, and topical administration. Preferably, the therapeutic polynucleotide composition contains an expression construct comprising a promoter operably linked to a polynucleotide of at least 12, 22, 25, 30, or 35 contiguous nt of the polynucleotide disclosed herein. Various methods can be used to administer the therapeutic composition directly to a specific site in the body. For example, a small metastatic lesion is located and the therapeutic composition injected several times in several different locations within the body of tumor. Alternatively, arteries which serve a tumor are identified, and the therapeutic composition injected into such an artery, in order to deliver the

composition directly into the tumor. A tumor that has a necrotic center is aspirated and the composition injected directly into the now empty center of the tumor. The antisense composition is directly administered to the surface of the tumor, for example, by topical application of the composition. X-ray imaging is used to assist in certain of the above delivery methods.

5 Receptor-mediated targeted delivery of therapeutic compositions containing an antisense polynucleotide, subgenomic polynucleotides, or antibodies to specific tissues can also be used. Receptor-mediated DNA delivery techniques are described in, for example, Findeis *et al.*, *Trends Biotechnol.* (1993) 11:202; Chiou *et al.*, *Gene Therapeutics: Methods And Applications Of Direct Gene Transfer* (J.A. Wolff, ed.) (1994); Wu *et al.*, *J. Biol. Chem.* (1988) 263:621; Wu *et al.*, *J. Biol. Chem.* (1994) 269:542; Zenke *et al.*, *Proc. Natl. Acad. Sci. (USA)* (1990) 87:3655; Wu *et al.*, *J. Biol. Chem.* (1991) 266:338. Therapeutic compositions containing a polynucleotide are administered in a range of about 100 ng to about 200 mg of DNA for local administration in a gene therapy protocol. Concentration ranges of about 500 ng to about 50 mg, about 1 g to about 2 mg, about 5 g to about 500 g, and about 20 g to about 100 g of DNA can also be used during a gene therapy protocol. Factors such as method of action (e.g., for enhancing or inhibiting levels of the encoded gene product) and efficacy of transformation and expression are considerations which will affect the dosage required for ultimate efficacy of the antisense subgenomic polynucleotides. Where greater expression is desired over a larger area of tissue, larger amounts of antisense subgenomic polynucleotides or the same amounts readministered in a successive protocol of administrations, or 15 several administrations to different adjacent or close tissue portions of, for example, a tumor site, may be required to effect a positive therapeutic outcome. In all cases, routine experimentation in clinical trials will determine specific ranges for optimal therapeutic effect. For polynucleotide-related genes encoding polypeptides or proteins with anti-inflammatory activity, suitable use, doses, and administration are described in USPN 5,654,173.

20 The therapeutic polynucleotides and polypeptides of the present invention can be delivered using gene delivery vehicles. The gene delivery vehicle can be of viral or non-viral origin (see generally, Jolly, *Cancer Gene Therapy* (1994) 1:51; Kimura, *Human Gene Therapy* (1994) 5:845; Connelly, *Human Gene Therapy* (1995) 1:185; and Kaplitt, *Nature Genetics* (1994) 6:148). Expression of such coding sequences can be induced using endogenous mammalian or heterologous 25 promoters. Expression of the coding sequence can be either constitutive or regulated.

30 Viral-based vectors for delivery of a desired polynucleotide and expression in a desired cell are well known in the art. Exemplary viral-based vehicles include, but are not limited to, recombinant retroviruses (see, e.g., WO 90/07936; WO 94/03622; WO 93/25698; WO 93/25234; USPN 5,219,740; WO 93/11230; WO 93/10218; USPN 4,777,127; GB Patent No. 2,200,651; EP 0 35 345 242; and WO 91/02805), alphavirus-based vectors (e.g., Sindbis virus vectors, Semliki forest

virus (ATCC VR-67: ATCC VR-1247). Ross River virus (ATCC VR-373: ATCC VR-1246) and Venezuelan equine encephalitis virus (ATCC VR-923: ATCC VR-1250: ATCC VR 1249: ATCC VR-532). and adeno-associated virus (AAV) vectors (see, e.g., WO 94/12649, WO 93/03769; WO 93/19191; WO 94/28938; WO 95/11984 and WO 95/00655). Administration of DNA linked to 5 killed adenovirus as described in Curiel, *Hum. Gene Ther.* (1992) 3:147 can also be employed.

Non-viral delivery vehicles and methods can also be employed, including, but not limited to, polycationic condensed DNA linked or unlinked to killed adenovirus alone (see, e.g., Curiel, *Hum. Gene Ther.* (1992) 3:147); ligand-linked DNA (see, e.g., Wu, *J. Biol. Chem.* (1989) 264:16985); eukaryotic cell delivery vehicles cells (see, e.g., USPN 5,814,482; WO 95/07994; WO 96/17072; 10 WO 95/30763; and WO 97/42338) and nucleic charge neutralization or fusion with cell membranes. Naked DNA can also be employed. Exemplary naked DNA introduction methods are described in WO 90/11092 and USPN 5,580,859. Liposomes that can act as gene delivery vehicles are described in USPN 5,422,120; WO 95/13796; WO 94/23697; WO 91/14445; and EP 0524968. Additional approaches are described in Philip, *Mol. Cell Biol.* (1994) 14:2411, and in Woffendin, *Proc. Natl. 15 Acad. Sci.* (1994) 91:1581

Further non-viral delivery suitable for use includes mechanical delivery systems such as the approach described in Woffendin *et al.*, *Proc. Natl. Acad. Sci. USA* (1994) 91(24):11581. Moreover, the coding sequence and the product of expression of such can be delivered through deposition of photopolymerized hydrogel materials or use of ionizing radiation (see, e.g., USPN 5,206,152 and 20 WO 92/11033). Other conventional methods for gene delivery that can be used for delivery of the coding sequence include, for example, use of hand-held gene transfer particle gun (see, e.g., USPN 5,149,655); use of ionizing radiation for activating transferred gene (see, e.g., USPN 5,206,152 and WO 92/11033).

The present invention will now be illustrated by reference to the following examples which 25 set forth particularly advantageous embodiments. However, it should be noted that these embodiments are illustrative and are not to be construed as restricting the invention in any way.

EXAMPLES

Example 1: Source of Biological Materials and Overview of Novel Polynucleotides Expressed 30 by the Biological Materials

cDNA libraries were constructed from either human colon cancer cell line Km12L4-A (Morikawa, *et al.*, *Cancer Research* (1988) 48:6863), KM12C (Morikawa *et al.* *Cancer Res.* (1988) 48:1943-1948), or MDA-MB-231 (Brinkley *et al.* *Cancer Res.* (1980) 40:3118-3129) was used to construct a cDNA library from mRNA isolated from the cells. Sequences expressed by these cell 35 lines were isolated and analyzed: most sequences were about 275-300 nucleotides in length. The

KM12L4-A cell line is derived from the KM12C cell line. The KM12C cell line, which is poorly metastatic (low metastatic) was established in culture from a Dukes' stage B₂ surgical specimen (Morikawa *et al.* *Cancer Res.* (1988) 48:6863). The KML4-A is a highly metastatic subline derived from KM12C (Yeatman *et al.* *Nucl. Acids. Res.* (1995) 23:4007; Bao-Ling *et al.* *Proc. Annu. Meet. Am. Assoc. Cancer. Res.* (1995) 21:3269). The KM12C and KM12C-derived cell lines (e.g., KM12L4, KM12L4-A, etc.) are well-recognized in the art as a model cell line for the study of colon cancer (see, e.g., Moriakawa *et al.*, *supra*; Radinsky *et al.* *Clin. Cancer Res.* (1995) 1:19; Yeatman *et al.*, (1995) *supra*; Yeatman *et al.* *Clin. Exp. Metastasis* (1996) 14:246). The MDA-MB-231 cell line was originally isolated from pleural effusions (Cailleau, *J. Natl. Cancer. Inst.* (1974) 53:661), is of high metastatic potential, and forms poorly differentiated adenocarcinoma grade II in nude mice consistent with breast carcinoma.

The sequences of the isolated polynucleotides were first masked to eliminate low complexity sequences using the XBLAST masking program (Claverie "Effective Large-Scale Sequence Similarity Searches." In: Computer Methods for Macromolecular Sequence Analysis, Doolittle, ed., *Meth. Enzymol.* 266:212-227 Academic Press, NY, NY (1996); see particularly Claverie, in "Automated DNA Sequencing and Analysis Techniques" Adams *et al.*, eds., Chap. 36, p. 267 Academic Press, San Diego, 1994 and Claverie *et al.* *Comput. Chem.* (1993) 17:191). Generally, masking does not influence the final search results, except to eliminate sequences of relative little interest due to their low complexity, and to eliminate multiple "hits" based on similarity to repetitive regions common to multiple sequences, e.g., Alu repeats. Masking resulted in the elimination of 43 sequences. The remaining sequences were then used in a BLASTN vs. GenBank search; sequences that exhibited greater than 70% overlap, 99% identity, and a p value of less than 1×10^{-40} were discarded. Sequences from this search also were discarded if the inclusive parameters were met, but the sequence was ribosomal or vector-derived.

The resulting sequences from the previous search were classified into three groups (1, 2 and 3 below) and searched in a BLASTX vs. NRP (non-redundant proteins) database search: (1) unknown (no hits in the GenBank search), (2) weak similarity (greater than 45% identity and p value of less than 1×10^{-5}), and (3) high similarity (greater than 60% overlap, greater than 80% identity, and p value less than 1×10^{-5}). Sequences having greater than 70% overlap, greater than 99% identity, and p value of less than 1×10^{-40} were discarded.

The remaining sequences were classified as unknown (no hits), weak similarity, and high similarity (parameters as above). Two searches were performed on these sequences. First, a BLAST vs. EST database search was performed and sequences with greater than 99% overlap,

greater than 99% similarity and a p value of less than 1×10^{-40} were discarded. Sequences with a p value of less than 1×10^{-65} when compared to a database sequence of human origin were also excluded. Second, a BLASTN vs. Patent GeneSeq database was performed and sequences having greater than 99% identity, p value less than 1×10^{-40} , and greater than 99% overlap were discarded.

5 The remaining sequences were subjected to screening using other rules and redundancies in the dataset. Sequences with a p value of less than 1×10^{-111} in relation to a database sequence of human origin were specifically excluded. The final result provided the 1,565 sequences listed as SEQ ID NOS:1-1565 in the accompanying Sequence Listing and summarized in Table 1A (inserted prior to claims). Each identified polynucleotide represents sequence from at least a partial mRNA transcript.

10 Table 1A provides: 1) the SEQ ID NO assigned to each sequence for use in the present specification; 2) the filing date of the U.S. priority application in which the sequence was first filed; 3) the attorney docket number assigned to the priority application (for internal use); 4) the SEQ ID NO assigned to the sequence in the priority application; 5) the sequence name used as an internal 15 identifier of the sequence; and 6) the name assigned to the clone from which the sequence was isolated. Because the provided polynucleotides represent partial mRNA transcripts, two or more polynucleotides of the invention may represent different regions of the same mRNA transcript and the same gene. Thus, if two or more SEQ ID NOS: are identified as belonging to the same clone, then either sequence can be used to obtain the full-length mRNA or gene.

15 In order to confirm the sequences of SEQ ID NOS:1-1565, the clones were retrieved from a library using a robotic retrieval system, and the inserts of the retrieved clones re-sequenced. These "validation" sequences are provided as SEQ ID NOS:1566-2610 in the Sequence Listing, and a summary of the "validation" sequences provided in Table 1B (inserted prior to claims). Table 1B provides: 1) the SEQ ID NO assigned to each sequence for use in the present specification; 2) the 20 sequence name assigned to the "validation" sequence obtained; 3) whether the "validation" sequence contains sequence that overlaps with an original sequence of SEQ ID NOS:1-1565 (Validation Overlap (VO)), or whether the "validation" sequence does not substantially overlap with an original sequence of SEQ ID NOS:1-1565 (indicated by Validation Non-Overlap (VNO)); and 25 4) where the sequence is indicated as VO, the name of the clone that contains the indicated "validation" sequence. "Validation" sequences are indicated as "VO" where the "validation" sequence overlaps with an original sequence (e.g., one of SEQ ID NOS:1-1565), and/or the "validation" sequence belongs to the same cluster as the original sequence using the clustering 30 technique described above. Because the inserts of the clones are generally longer than the original

sequence and the validation sequence, it is possible that a "validation" sequence can be obtained from the same clone as an original sequence but yet not share any of the sequence of the original. Such validation sequences will, however, belong to the same cluster as the original sequence using the clustering technique described above. VO "validation" sequences are contained within the same 5 clone as the original sequence (one of SEQ ID NOS:1-1565). "Validation" sequences that provided overlapping sequence are indicating by "VO" can be correlated with the original sequences they validate by referring to Table 1A. Sequences indicated as VNO are treated as newly isolated sequences and may or may not be related to the sequences of SEQ ID NOS:1-1565. Because the 10 "validation" sequences are often longer than the original polynucleotide sequences and thus provide additional sequence information. All validation sequences can be obtained either from an indicated clone (e.g., for VO sequences) or from a cDNA library described herein (e.g., using primers designed from the sequence provided in the sequence listing).

Example 2: Results of Public Database Search to Identify Function of Gene Products

15 SEQ ID NOS:1566-2610 were translated in all three reading frames, and the nucleotide sequences and translated amino acid sequences used as query sequences to search for homologous sequences in either the GenBank (nucleotide sequences) or Non-Redundant Protein (amino acid sequences) databases. Query and individual sequences were aligned using the BLAST 2.0 programs, available over the world wide web at <http://www.ncbi.nlm.nih.gov/BLAST/>. (see also Altschul, et al. 20 *Nucleic Acids Res.* (1997) 25:3389-3402). The sequences were masked to various extents to prevent searching of repetitive sequences or poly-A sequences, using the XBLAST program for masking low complexity as described above in Example 1.

Tables 2A and 2B (inserted before the claims) provide the alignment summaries having a p value of 1×10^{-2} or less indicating substantial homology between the sequences of the present 25 invention and those of the indicated public databases. Table 2A provides the SEQ ID NO of the query sequence, the accession number of the GenBank database entry of the homologous sequence, and the p value of the alignment. Table 2A provides the SEQ ID NO of the query sequence, the accession number of the Non-Redundant Protein database entry of the homologous sequence, and the p value of the alignment. The alignments provided in Tables 2A and 2B are the best available 30 alignment to a DNA or amino acid sequence at a time just prior to filing of the present specification. The activity of the polypeptide encoded by the SEQ ID NOS listed in Tables 2A and 2B can be extrapolated to be substantially the same or substantially similar to the activity of the reported nearest neighbor or closely related sequence. The accession number of the nearest neighbor is reported, providing a publicly available reference to the activities and functions exhibited by the

nearest neighbor. The public information regarding the activities and functions of each of the nearest neighbor sequences is incorporated by reference in this application. Also incorporated by reference is all publicly available information regarding the sequence, as well as the putative and actual activities and functions of the nearest neighbor sequences listed in Table 2 and their related sequences. The search program and database used for the alignment, as well as the calculation of the p value are also indicated.

Full length sequences or fragments of the polynucleotide sequences of the nearest neighbors can be used as probes and primers to identify and isolate the full length sequence of the corresponding polynucleotide. The nearest neighbors can indicate a tissue or cell type to be used to construct a library for the full-length sequences of the corresponding polynucleotides.

Example 3: Members of Protein Families

SEQ ID NOS:1566-2601 were used to conduct a profile search as described in the specification above. Several of the polynucleotides of the invention were found to encode polypeptides having characteristics of a polypeptide belonging to a known protein family (and thus represent new members of these protein families) and/or comprising a known functional domain (Table 3A, inserted prior to claims). Table 3A provides the SEQ ID NO: of the query sequence, a brief description of the profile hit, the position of the query sequence within the individual sequence (indicated as "start" and "stop"), and the orientation (Direction) of the query sequence with respect to the individual sequence. where forward (for) indicates that the alignment is in the same direction (left to right) as the sequence provided in the Sequence Listing and reverse (rev) indicates that the alignment is with a sequence complementary to the sequence provided in the Sequence Listing.

Some polynucleotides exhibited multiple profile hits where the query sequence contains overlapping profile regions, and/or where the sequence contains two different functional domains. Each of the profile hits of Table 3A are described in more detail below. The acronyms for the profiles (provided in parentheses) are those used to identify the profile in the Pfam and Prosite databases. The Pfam database can be accessed through any of the following URLs: <http://pfam.wustl.edu/index.html>; <http://www.sanger.ac.uk/Software/Pfam/>; and <http://www.cgr.ki.se/Pfam/>. The Prosite database can be accessed at <http://www.expasy.ch/prosite/>. The public information available on the Pfam and Prosite databases regarding the various profiles, including but not limited to the activities, function, and consensus sequences of various protein families and protein domains, is incorporated herein by reference.

14-3-3 Family (14_3_3), SEQ ID NO:1967 corresponds to a sequence encoding a 14-3-3 protein family member. The 14-3-3 protein family includes a group of closely related acidic homodimeric proteins of about 30 kD first identified as very abundant in mammalian brain tissues

and located preferentially in neurons (Aitken et al. *Trends Biochem. Sci.* (1995) 20:95-97; Morrison *Science* (1994) 266:56-57; and Xiao et al. *Nature* (1995) 376:188-191). The 14-3-3 proteins have multiple biological activities, including a key role in signal transduction pathways and the cell cycle.

14-3-3 proteins interact with kinases (e.g., PKC or Raf-1), and can also function as protein-kinase dependent activators of tyrosine and tryptophan hydroxylases. The 14-3-3 protein sequences are extremely well conserved, and include two highly conserved regions: the first is a peptide of 11 residues located in the N-terminal section; the second, a 20 amino acid region located in the C-terminal section. The consensus patterns are as follows: 1) R-N-L-[LIV]-S-[VG]-[GA]-Y-[KN]-N-[IVA]; 2) Y-K-[DE]-S-T-L-I-[IM]-Q-L-[LF]-[RHC]-D-N-[LF]-T-[LS]-W-[TAN]-[SAD].

10 **3'5'-Cyclin Nucleotide Phosphodiesterases (PDEase)**. SEQ ID NO: 2366 represents a polynucleotide encoding a novel 3'5'-cyclic nucleotide phosphodiesterase. PDEases catalyze the hydrolysis of cAMP or cGMP to the corresponding nucleoside 5' monophosphates (Charbonneau et al. *Proc. Natl. Acad. Sci. U.S.A.* (1986) 83:9308). There are at least seven different subfamilies of PDEases (Beavo et al., *Trends Pharmacol. Sci.* (1990) 11:150; <http://weber.u.washington.edu/~pde/>):

15 1) Type 1, calmodulin/calcium-dependent PDEases; 2) Type 2, cGMP-stimulated PDEases; 3) Type 3, cGMP-inhibited PDEases; 4) Type 4, cAMP-specific PDEases; 5) Type 5, cGMP-specific PDEases; 6) Type 6, rhodopsin-sensitive cGMP-specific PDEases; and 7) Type 7, High affinity cAMP-specific PDEases. All PDEase forms share a conserved domain of about 270 residues. The signature pattern is determined from a stretch of 12 residues that contains two conserved histidines:

20 H-D-[LIVMFY]-x-H-x-[AG]-x(2)-[NQ]-x-[LIVMFY].

25 **Four Transmembrane Integral Membrane Proteins (transmembrane4)**. SEQ ID NOS:1579 and 1978 sequences correspond to a sequence encoding a member of the four transmembrane segments integral membrane protein family (tm4 family). The tm4 family of proteins includes a number of evolutionarily-related eukaryotic cell surface antigens (Levy et al., *J. Biol. Chem.* (1991) 266:14597; Tomlinson et al., *Eur. J. Immunol.* (1993) 23:136; Barclay et al. The leucocyte antigen factbooks. (1993) Academic Press, London/San Diego). The tm4 family members are type III membrane proteins, which are integral membrane proteins containing an N-terminal membrane-anchoring domain that functions both as a translocation signal and as a membrane anchor. The family members also contain three additional transmembrane regions, at least seven conserved 30 cysteines residues, and are of approximately the same size (218 to 284 residues). The consensus pattern spans a conserved region including two cysteines located in a short cytoplasmic loop between two transmembrane domains: Consensus pattern: G-x(3)-[LIVMF]-x(2)-[GSA]-[LIVMF](2)-G-C-x-[GA]-[STA]-x(2)-[EG]-x(2)-[CWN]-[LIVM](2).

35 **Seven Transmembrane Integral Membrane Proteins -- Rhodopsin Family (7tm_1)**. SEQ ID NOS:1652, 1927, and 2068 correspond to a sequence encoding a member of the seven

transmembrane (7tm) receptor rhodopsin family. G-protein coupled receptors of the (7tm) rhodopsin family include hormones, neurotransmitters, and light receptors that transduce extracellular signals by interaction with guanine nucleotide-binding (G) proteins (Strosberg *Eur. J. Biochem.* (1991) 196:1, Kerlavage *Curr. Opin. Struct. Biol.* (1991) 1:394. Probst, et al., *DNA Cell Biol.* (1992) 11:1, Savarese, et al., *Biochem. J.* (1992) 283:1, <http://www.ecrdb.uthscsa.edu/>, <http://swiss.embl-heidelberg.de/7tm/>) The consensus pattern that contains the conserved triplet and that also spans the major part of the third transmembrane helix is used to detect this widespread family of proteins: [GSTALIVMFYWC]-[GSTANCPDE]-{EDPKRH}-x(2)-[LIVMNQGA]-x(2)-[LIVMFT]-[GSTANC]-[LIVMFYWSTAC]-[DENH]-R-[FYWCSH]-x(2)-[LIVM].

10 Seven Transmembrane Integral Membrane Proteins -- Secretin Family (7tm 2), SEQ ID NOS:1598, 1719, 1911, 1927, 2068, and 2341 correspond to a sequence encoding a member of the seven transmembrane receptor (7tm) secretin family (Jueppner et al. *Science* (1991) 254:1024; Hamann et al. *Genomics* (1996) 32:144). The N-terminal extracellular domain of these receptors contains five conserved cysteines residues involved in disulfide bonds, with a consensus pattern in 15 the region that spans the first three cysteines. One of the most highly conserved regions spans the C-terminal part of the last transmembrane region and the beginning of the adjacent intracellular region and is used as a second signature pattern. The two consensus patterns are: 1) C-x(3)-[FYWLIV]-D-x(3,4)-C-[FW]-x(2)-[STAGV]-x(8,9)-C-[PF]; and 2) Q-G-[LMFCA]-[LIVMFT]-[LIV]-x-[LIVFST]-[LIF]-[VFYH]-C-[LFY]-x-N-x(2)-V

20 ATPases Associated with Various Cellular Activities (ATPases). Several of the polynucleotides of the invention correspond to a sequence that encodes a member of a family of ATPases Associated with diverse cellular Activities (AAA). The AAA protein family is composed of a large number of ATPases that share a conserved region of about 220 amino acids containing an ATP-binding site (Froehlich et al., *J. Cell Biol.* (1991) 114:443; Erdmann et al. *Cell* (1991) 64:499; Peters et al., *EMBO J.* (1990) 9:1757; Kunau et al., *Biochimie* (1993) 75:209-224; Confalonieri et al., *BioEssays* (1995) 17:639; <http://yeamob.pci.chemie.uni-tuebingen.de/AAA/Description.html>). The AAA domain, which can be present in one or two copies, acts as an ATP-dependent protein clamp (Confalonieri et al. (1995) *BioEssays* 17:639) and contains a highly conserved region located in the central part of the domain. The consensus pattern is: [LIVMT]-x-[LIVMT]-[LIVMF]-x-[GATMC]-[ST]-[NS]-x(4)-[LIVM]-D-x-A-[LIFA]-x-R.

30 Basic Region Plus Leucine Zipper Transcription Factors (BZIP), SEQ ID NO:1623 represents a polynucleotide encoding a novel member of the family of basic region plus leucine zipper transcription factors. The bZIP superfamily (Hurst, *Protein Prof.* (1995) 2:105; and Ellenberger, *Curr. Opin. Struct. Biol.* (1994) 4:12) of eukaryotic DNA-binding transcription factors 35 encompasses proteins that contain a basic region mediating sequence-specific DNA-binding

followed by a leucine zipper required for dimerization. The consensus pattern for this protein family is: [KR]-x(1,3)-[RKSAQ]-N-x(2)-[SAQ](2)-x-[RKTAENQ]-x-R-x-[RK].

C2 domain (C2). SEQ ID NOS: 1715 and 2426 correspond to a sequence encoding a C2 domain, which is involved in calcium-dependent phospholipid binding (Davletov *J. Biol. Chem.* 1993) 268:26386-26390) or, in proteins that do not bind calcium, the domain may facilitate binding to inositol-1,3,4,5-tetraphosphate (Fukuda et al. *J. Biol. Chem.* (1994) 269:29206-29211; Sutton et al. *Cell* (1995) 80:929-938). The consensus sequence is: [ACG]-x(2)-L-x(2,3)-D-x(1,2)-[NGSTLIF]-[GTMR]-x-[STAP]-D- [PA]-[FY].

Cysteine proteases (Cys-protease). SEQ ID NO:2238 represents a polynucleotide encoding a protein having a eukaryotic thiol (cysteine) protease active site. Cysteine proteases (Dufour *Biochimie* (1988) 70:1335) are a family of proteolytic enzymes that contain an active site cysteine. Catalysis proceeds through a thioester intermediate and is facilitated by a nearby histidine side chain: an asparagine completes the essential catalytic triad. The sequences around the three active site residues are well conserved and can be used as signature patterns: Q-x(3)-[GE]-x-C-[YW]-x(2)-[STAGC]-[STAGCV] (where C is the active site residue); 2) [LIVMGSTAN]-x-H-[GSACE]-[LIVM]-x-[LIVMAT](2)-G-x-[GSADNH] (where H is the active site residue); and 3) [FYCH]-[WI]-[LIVT]-x-[KIQAG]-N-[ST]-W-x(3)-[FYW]-G-x(2)-G- [LFYW]-[LIVMFY]-x-[LIVMF] (where N is the active site residue).

DEAD and DEAH box families ATP-dependent helicases (Dead_box_helic). SEQ ID NOS:1630, 1865, and 2517 represent polynucleotides encoding a novel member of the DEAD and DEAH box families (Schmid et al., *Mol. Microbiol.* (1992) 6:283; Linder et al., *Nature* (1989) 337:121; Wasserman, et al., *Nature* (1991) 349:463). All members of these families are involved in ATP-dependent, nucleic-acid unwinding. All DEAD box family members share a number of conserved sequence motifs, some of which are specific to the DEAD family, with others shared by other ATP-binding proteins or by proteins belonging to the helicases 'superfamily' (Hodgman *Nature* (1988) 333:22 and *Nature* (1988) 333:578 (Errata); http://www.expasy.ch/www/linder/HELICASES_TEXT.html). One of these motifs, called the 'D-E-A-D-box', represents a special version of the B motif of ATP-binding proteins. Proteins that have His instead of the second Asp and are 'D-E-A-H-box' proteins (Wasserman et al., *Nature* (1991) 349:463; Harosh, et al., *Nucleic Acids Res.* (1991) 19:6331; Koonin , et al., *J. Gen. Virol.* (1992) 73:989; http://www.expasy.ch/www/linder/HELICASES_TEXT.html). The following signature patterns are used to identify member for both subfamilies: 1) [LIVMF](2)-D-E-A-D-[RKEN]-x-[LIVMFYGSTN]; and 2) [GSAH]-x-[LIVMF](3)-D-E-[ALIV]-H-[NECR].

Dual specificity phosphatase (DSPc). Dual specificity phosphatases (DSPs) are Ser/Thr and Tyr protein phosphatases that comprise a tertiary fold highly similar to that of tyrosine-specific

phosphatases, except for a "recognition" region connecting helix alpha1 to strand beta1. This tertiary fold may determine differences in substrate specific between VH-I related dual specificity phosphatase (VHR), the protein tyrosine phosphatases (PTPs), and other DSPs. Phosphatases are important in the control of cell growth, proliferation, differentiation and transformation.

5 EF Hand (EFhand). SEQ ID NO:1595 corresponds to a polynucleotide encoding a member of the EF-hand protein family, a calcium binding domain shared by many calcium-binding proteins belonging to the same evolutionary family (Kawasaki *et al.* *Protein. Prof.* (1995) 2:305-490). The domain is a twelve residue loop flanked on both sides by a twelve residue alpha-helical domain, with a calcium ion coordinated in a pentagonal bipyramidal configuration. The six residues involved in 10 the binding are in positions 1, 3, 5, 7, 9 and 12; these residues are denoted by X, Y, Z, -Y, -X and -Z. The invariant Glu or Asp at position 12 provides two oxygens for liganding Ca (bidentate ligand). The consensus pattern includes the complete EF-hand loop as well as the first residue which follows the loop and which seem to always be hydrophobic: D-x-[DNS]-{ILVFYW}-{DENSTG]-[DNQGHRK]-{GP}-{LIVMC}-{DENQSTAGC}-x(2)-[DE]-[LIVMFYW].

15 Eukaryotic Aspartyl Proteases (asp). Several of the polynucleotides of the invention correspond to a sequence encoding a novel eukaryotic aspartyl protease. Aspartyl proteases, known as acid proteases, (EC 3.4.23.-) are a widely distributed family of proteolytic enzymes (Foltmann., *Essays Biochem.* (1981) 17:52; Davies. *Annu. Rev. Biophys. Chem.* (1990) 19:189; Rao, *et al.*, *Biochemistry* (1991) 30:4663) known to exist in vertebrates, fungi, plants, retroviruses and some 20 plant viruses. Aspartate proteases of eukaryotes are monomeric enzymes which consist of two domains. Each domain contains an active site centered on a catalytic aspartyl residue. The consensus pattern to identify eukaryotic aspartyl protease is: [LIVMFGAC]-[LIVMTADN]-[LIVFSA]-D-[ST]-G-[STAV]-[STAPDENQ]-x-[LIVMFSTNC]-x-[LIVMFGTA], where D is the active site residue.

25 Fibronectin Type II collagen-binding domain (FntypeII). SEQ ID NO: 1968 corresponds to a polynucleotide encoding a polypeptide having a type II fibronectin collagen binding domain. Fibronectin is a plasma protein that binds cell surfaces and various compounds including collagen, fibrin, heparin, DNA, and actin. The major part of the sequence of fibronectin consists of the repetition of three types of domains, called type I, II, and III (Skorstengaard *et al.*, *Eur. J. Biochem.* 30 (1986) 161:441). The type II domain, which is duplicated in fibronectin, is approximately forty residues long, contains four conserved cysteines involved in disulfide bonds and is part of the collagen-binding region of fibronectin. The consensus pattern for identifying members of this family, which pattern spans this entire domain, is: C-x(2)-P-F-x-[FYWI]-x(7)-C-x(8,10)-W-C-x(4)-[DNSR]-[FYW]-x(3,5)-[FYW]-x-[FYWI]-C (where the four C's are involved in disulfide bonds).

35 G-Protein Alpha Subunit (G-alpha). SEQ ID NO: 1779 corresponds to a gene encoding a

member of the G-protein alpha subunit family. G-proteins are a family of membrane-associated proteins that couple extracellularly-activated integral-membrane receptors to intracellular effectors, such as ion channels and enzymes that vary the concentration of second messenger molecules. G-proteins are composed of 3 subunits (alpha, beta and gamma) which, in the resting state, associate as a trimer at the inner face of the plasma membrane. The alpha subunit, which binds GTP and exhibits GTPase activity, is about 350-400 amino acids in length with a molecular weight in the range of 40-45 kDa. Seventeen distinct types of alpha subunit have been identified in mammals, and fall into 4 main groups on the basis of both sequence similarity and function: alpha-s, alpha-q, alpha-i and alpha-12 (Simon *et al.*, *Science* (1993) 252:802). They are often N-terminally acylated, usually with myristate and/or palmitoylate, and these fatty acid modifications can be important for membrane association and high-affinity interactions with other proteins.

10 Helicases conserved C-terminal domain (helicase_C). SEQ ID NOS: 1621 and 1652 represent polynucleotides encoding novel members of the DEAD/H helicase family. The DEAD and DEAH families are described above.

15 Helix-Loop-Helix (HLH) DNA Binding Domain (HLH). SEQ ID NO:2192 corresponds to a sequence encoding an HLH domain. The HLH domain, which normally spans about 40 to 50 amino acids, is present in a number of eukaryotic transcription factors. The HLH domain is formed of two amphipathic helices joined by a variable length linker region that forms a loop that mediates protein dimerization (Murre *et al.*, *Cell* (1989) 56:777-783). Basic HLH proteins (bHLH), which 20 have an extra basic region of about 15 amino acid residues adjacent the HLH domain and specifically bind to DNA, include two groups: class A (ubiquitous) and class B (tissue-specific). bHLH family members bind variations of the E-box motif (CANNTG). The homo- or heterodimerization mediated by the HLH domain is independent of, but necessary for DNA binding, as two basic regions are required for DNA binding activity. The HLH proteins lacking the basic 25 domain function as negative regulators since they form heterodimers, but fail to bind DNA. Consensus pattern: [DENSTAP]-[KTR]-[LIVMAGSNT]-{FYWCPHKR}-[LIVMT]-[LIVM]-x(2)-[STAV]-[LIVMSTACKR]-x-[VMFYH]-[LIVMTA]-{P}-{P}-[LIVMRKHQ].

30 Kinase Domain of Tors. The TOR profile is directed towards a lipid kinase protein family. This family is composed of large proteins with a lipid and protein kinase domain and characterized through their sensitivity to rapamycin (an antifungal compound). TOR proteins are involved in signal transduction downstream of PI3 kinase and many other signals. TOR (also called FRAP, RAFT) plays a role in regulating protein synthesis and cell growth, and in yeast controls translation initiation and early G1 progression. See, e.g., Barbet *et al.* *Mol Biol Cell*. (1996) 7(1):25-42; Helliwell *et al.* *Genetics* (1998) 148:99-112.

35 MAP kinase kinase (mkk). SEQ ID NOS: 1825, 1876, 2039, and 2526 represent members of

the MAP kinase kinase (mkk) family. MAP kinases (MAPK) are involved in signal transduction, and are important in cell cycle and cell growth controls. The MAP kinase kinases (MAPKK) are dual-specificity protein kinases which phosphorylate and activate MAP kinases. MAPKK homologues have been found in yeast, invertebrates, amphibians, and mammals. Moreover, the 5 MAPKK/MAPK phosphorylation switch constitutes a basic module activated in distinct pathways in yeast and in vertebrates. MAPKKs are essential transducers through which signals must pass before reaching the nucleus. For review, see, e.g., Biologique *Biol Cell* (1993) 79:193-207; Nishida *et al.*, *Trends Biochem Sci* (1993) 18:128-31; Ruderman *Curr Opin Cell Biol* (1993) 5:207-13; Dhanasekaran *et al.*, *Oncogene* (1998) 17:1447-55; Kiefer *et al.*, *Biochem Soc Trans* (1997) 25:491-10 8; and Hill, *Cell Signal* (1996) 8:533-44.

15 Neurotransmitter-Gated Ion-Channel (neur_chan). Several of the sequences correspond to a sequence encoding a neurotransmitter-gated ion channel. Neurotransmitter-gated ion-channels, which provide the molecular basis for rapid signal transmission at chemical synapses, are post-synaptic oligomeric transmembrane complexes that transiently form a ionic channel upon the binding of a specific neurotransmitter. Five types of neurotransmitter-gated receptors are known: 1) nicotinic acetylcholine receptor (AchR); 2) glycine receptor; 3) gamma-aminobutyric-acid (GABA) receptor; 4) serotonin 5HT3 receptor; and 5) glutamate receptor. All known sequences of subunits from neurotransmitter-gated ion-channels are structurally related, and are composed of a large extracellular glycosylated N-terminal ligand-binding domain, followed by three hydrophobic transmembrane regions that form the ionic channel, followed by an intracellular region of variable length. A fourth hydrophobic region is found at the C-terminal of the sequence. The consensus pattern is: C-x-[LIVMFQ]-x-[LIVMF]-x(2)-[FY]-P-x-D-x(3)-C, where the two C's are linked by a disulfide bond.

Protein Kinase (protkinase). Several sequences represent polynucleotides encoding protein kinases, which catalyze phosphorylation of proteins in a variety of pathways, and are implicated in cancer. Eukaryotic protein kinases (Hanks, *et al.*, *FASEB J.* (1995) 9:576; Hunter, *Meth. Enzymol.* (1991) 200:3; Hanks, *et al.*, *Meth. Enzymol.* (1991) 200:38; Hanks, *Curr. Opin. Struct. Biol.* (1991) 1:369; Hanks *et al.*, *Science* (1988) 241:42) belong to a very extensive family of proteins that share a conserved catalytic core common to both serine/threonine and tyrosine protein kinases. There are a number of conserved regions in the catalytic domain of protein kinases. The first region, located in the N-terminal extremity of the catalytic domain, is a glycine-rich stretch of residues in the vicinity of a lysine residue, which has been shown to be involved in ATP binding. The second region, located in the central part of the catalytic domain, contains a conserved an aspartic acid residue that is important for the catalytic activity of the enzyme (Knighton, *et al.*, *Science* (1991) 253:407).

35 The protein kinase profile includes two signature patterns for this second region: one

specific for serine/threonine kinases and the other for tyrosine kinases. A third profile is based on the alignment in (Hanks, *et al.*, *FASEB J.* (1995) 9:576) and covers the entire catalytic domain. The consensus patterns are as follows: 1) [LIV]-G-{P}-G-{P}-[FYWMGSTNH]-[SGA]-{PW}-[LIVCAT]-{PD}-x-[GSTACLIVMFY]-x(5.18)-[LIVMFYWCSTAR]-[AIVP]-[LIVMFAGCKR]-K, 5 where K binds ATP; 2) [LIVMFYC]-x-[HY]-x-D-[LIVMFY]-K-x(2)-N-[LIVMFYCT](3), where D is an active site residue; and 3) [LIVMFYC]-x-[HY]-x-D-[LIVMFY]-[RSTAC]-x(2)-N-[LIVMFYC], where D is an active site residue.

Protein Tyrosine Phosphatase (Y phosphatase) (PTPase). SEQ ID NOS: 1719, 1769, 2062, 2197, and 2275 represent polynucleotides encoding a tyrosine-specific protein phosphatase, a kinase 10 that catalyzes the removal of a phosphate groups attached to a tyrosine residue (EC 3.1.3.48) (PTPase) (Fischer *et al.*, *Science* (1991) 253:401; Charbonneau *et al.*, *Annu. Rev. Cell Biol.* (1992) 8:463; Trowbridge *Biol. Chem.* (1991) 266:23517; Tonks *et al.*, *Trends Biochem. Sci.* (1989) 14:497; and Hunter, *Cell* (1989) 58:1013). PTPases are important in the control of cell growth, proliferation, differentiation and transformation. Multiple forms of PTPase have been characterized and can be 15 classified into two categories: soluble PTPases and transmembrane receptor proteins that contain PTPase domain(s). Structurally, all known receptor PTPases are made up of a variable length extracellular domain, followed by a transmembrane region and a C-terminal catalytic cytoplasmic domain. PTPase domains consist of about 300 amino acids. Two conserved cysteines are absolutely required for activity, with a number of other conserved residues in the immediate vicinity also 20 important for activity. The consensus pattern for PTPases is: [LIVMF]-H-C-x(2)-G-x(3)-[STC]-[STAGP]-x-[LIVMFY]; C is the active site residue.

RNA Recognition Motif (RRM). SEQ ID NOS: 1850 and 2194 correspond to sequence 25 encoding an RNA recognition motif, also known as an RRM, RBD, or RNP domain. This domain, which is about 90 amino acids long, is contained in eukaryotic proteins that bind single-stranded RNA (Bandziulis *et al.*, *Genes Dev.* (1989) 3:431-437; Dreyfuss *et al.*, *Trends Biochem. Sci.* (1988) 13:86-91). Two regions within the RNA-binding domain are highly conserved: the first is a hydrophobic segment of six residues (which is called the RNP-2 motif), the second is an octapeptide motif (which is called RNP-1 or RNP-CS). The consensus pattern is: [RK]-G-{EDRKHPCG}-[AGSCI]-[FY]-[LIVA]-x-[FYLM].

SH2 Domain (SH2). SEQ ID NO: 2441 corresponds to a sequence encoding an SH2 30 domain. The Src homology 2 (SH2) domain includes an approximately 100 amino acid residue domain, which is conserved in the oncoproteins Src and Fps, as well as in many other intracellular signal-transducing proteins (Sadowski *et al.*, *Mol. Cell. Biol.* (1986) 6:4396-4408; Russel *et al.*, *FEBS Lett.* (1992) 304:15-20). SH2 domains function as regulatory modules of intracellular 35 signaling cascades by interacting with high affinity to phosphotyrosine-containing target peptides in

a sequence-specific and strictly phosphorylation-dependent manner. The SH2 domain has a conserved 3D structure consisting of two alpha helices and six to seven beta-strands. The core of the domain is formed by a continuous beta-meander composed of two connected beta-sheets (Kuriyan et al. *Curr. Opin. Struct. Biol.* (1993) 3:828-837).

5 Thioredoxin family active site (Thioredox). SEQ ID NO: 1618 represents a polynucleotide encoding a protein of the thioredoxin family. Thioredoxins are small proteins of approximately one hundred amino acid residues that participate in various redox reactions via the reversible oxidation of an active center disulfide bond (Holmgren. *Annu. Rev. Biochem.* (1985) 54:237; Gleason, et al., *FEMS Microbiol. Rev.* (1988) 54:271; Holmgren A. *J. Biol. Chem.* (1989) 264:13963; Eklund, et al.

10 *Proteins* (1991) 11:13). Thioredoxins exist in either reduced or oxidized forms where the two cysteine residues are linked in an intramolecular disulfide bond. The sequence around the redox-active disulfide bond is well conserved. The consensus pattern is: [LIVMF]-[LIVMSTA]-x-[LIVMFYC]-[FYWSTHE]-x(2)-[FYWGTN]-C- [GATPLVE]-[PHYWSTA]-C-x(6)-[LIVMFYWT] (where the two C's form the redox-active bond).

15 Trypsin (trypsin). SEQ ID NOS: 1579, 2290, 2341, 2421, 2430, and 2438 correspond to novel serine proteases of the trypsin family. The catalytic activity of the serine proteases from the trypsin family is provided by a charge relay system involving an aspartic acid residue hydrogen-bonded to a histidine, which itself is hydrogen-bonded to a serine. The sequences in the vicinity of the active site serine and histidine residues are well conserved (Brenner *Nature* (1988) 334:528).

20 The consensus patterns for the trypsin protein family are: 1) [LIVM]-[ST]-A-[STAG]-H-C, where H is the active site residue; and 2) [DNSTAGC]-[GSTAPIMVQH]-x(2)-G-[DE]-S-G-[GS]-[SAPHV]-[LIVMFYWH]-[LIVMFYSTANQH], where S is the active site residue. All sequences known to belong to this family are detected by the above consensus sequences, except for 18 different proteases which have lost the first conserved glycine. If a protein includes both the serine and the histidine active site signatures, the probability of it being a trypsin family serine protease is 100%.

25

30 WD Domain, G-Beta Repeats (WD_domain). SEQ ID NO: 2281 represents a members of the WD domain/G-beta repeat family. Beta-transducin (G-beta) is one of the three subunits (alpha, beta, and gamma) of the guanine nucleotide-binding proteins (G proteins) which act as intermediaries in the transduction of signals generated by transmembrane receptors (Gilman, *Annu. Rev. Biochem.* (1987) 56:615). The alpha subunit binds to and hydrolyzes GTP; the beta and gamma subunits are required for the replacement of GDP by GTP as well as for membrane anchoring and receptor recognition. In higher eukaryotes, G-beta exists as a small multigene family of highly conserved proteins of about 340 amino acid residues. Structurally, G-beta has eight tandem repeats of about 40 residues, each containing a central Trp-Asp motif (this type of repeat is sometimes called a WD-40 repeat). The consensus pattern for the WD domain/G-Beta repeat family is:

[LIVMSTAC]-[LIVMFYWSTAGC]-[LIMSTAG]-[LIVMSTAGC]-x(2)-[DN]-x(2)-
[LIVMWSTAC]-x-[LIVMFSTAG]-W-[DEN]-[LIVMFSTAGCN].

wnt Family of Developmental Signaling Proteins (Wnt_dev_sign). Several of the sequences correspond to novel members of the wnt family of developmental signaling proteins. Wnt-1 (previously known as int-1), the seminal member of this family, (Nusse. *Trends Genet.* (1988) 4:291) plays a role in intercellular communication and is important in central nervous system development. All wnt family proteins share the following features characteristic of secretory proteins: a signal peptide, several potential N-glycosylation sites and 22 conserved cysteines that may be involved in disulfide bonds. Wnt proteins generally adhere to the plasma membrane of secreting cells and are therefore likely to signal over only few cell diameters. The consensus pattern, which is based upon a highly conserved region including three cysteines, is as follows: C-K-
C-H-G-[LIVMT]-S-G-x-C.

Zinc Finger, C2H2 Type (Zincfing_C2H2). SEQ ID NOS: 1735, 1942, 2018, 2254, and 2515 correspond to polynucleotides encoding members of the C2H2 type zinc finger protein family, which contain zinc finger domains that facilitate nucleic acid binding (Klug *et al.*, *Trends Biochem. Sci.* (1987) 12:464; Evans *et al.*, *Cell* (1988) 52:1; Payre *et al.*, *FEBS Lett.* (1988) 234:245; Miller *et al.*, *EMBO J.* (1985) 4:1609; and Berg, *Proc. Natl. Acad. Sci. USA* (1988) 85:99). In addition to the conserved zinc ligand residues, a number of other positions are also important for the structural integrity of the C2H2 zinc fingers. (Rosenfeld *et al.*, *J. Biomol. Struct. Dyn.* (1993) 11:557) The best conserved position, which is generally an aromatic or aliphatic residue, is located four residues after the second cysteine. The consensus pattern for C2H2 zinc fingers is: C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H. The two C's and two H's are zinc ligands.

Example 4: Differential Expression of Polynucleotides of the Invention: Description of Libraries and Detection of Differential Expression

The relative expression levels of the polynucleotides of the invention was assessed in several libraries prepared from various sources, including cell lines and patient tissue samples. Table 4 provides a summary of these libraries, including the shortened library name (used hereafter), the mRNA source used to prepared the cDNA library, the "nickname" of the library that is used in the tables below (in quotes), and the approximate number of clones in the library.

Table 4. Description of cDNA Libraries

Library (lib #)	Description	Number of Clones in Cluster
1	Km12 L4 Human Colon Cell Line, High Metastatic Potential (derived from Km12C): "High Met Colon"	307133

Library (lib #)	Description	Number of Cl lines in Cluster
2	KM12C Human Colon Cell Line. Low Metastatic Potential; "Low Met Colon"	284755
3	MDA-MB-231 Human Breast Cancer Cell Line. High Metastatic Potential: micro-metastases in lung: "High Met Breast"	326937
4	MCF7 Human Breast Cancer Cell. Non Metastatic: "Low Met Breast"	318979
8	MV-522 Human Lung Cancer Cell Line. High Metastatic Potential: "High Met Lung"	223620
9	UCP-3 Human Lung Cancer Cell Line. Low Metastatic Potential: "Low Met Lung"	312503
12	Human microvascular endothelial cells (HMEC) – Untreated PCR (OligodT) cDNA library; "HMEC"	41938
13	Human microvascular endothelial cells (HMEC) – Basic fibroblast growth factor (bFGF) treated PCR (OligodT) cDNA library: "HMEC-bFGF"	42100
14	Human microvascular endothelial cells (HMEC) – Vascular endothelial growth factor (VEGF) treated PCR (OligodT) cDNA library; "HMEC-VEGF"	42825
15	Normal Colon – UC#2 Patient PCR (OligodT) cDNA library; "Normal Colon Tissue"	282722
16	Colon Tumor – UC#2 Patient PCR (OligodT) cDNA library; "Normal Colon Tumor Tissue"	298831
17	Liver Metastasis from Colon Tumor of UC#2 Patient PCR (OligodT) cDNA library; "High Met Colon Tissue"	303467
18	Normal Colon – UC#3 Patient PCR (OligodT) cDNA library; "Normal Colon Tissue"	36216
19	Colon Tumor – UC#3 Patient PCR (OligodT) cDNA library; "Colon Tumor Tissue"	41388
20	Liver Metastasis from Colon Tumor of UC#3 Patient PCR (OligodT) cDNA library; "High Met Colon Tissue"	30956
21	GRRpz Human Prostate Cell Line: "Normal Prostate"	164801
22	Woca Human Prostate Cancer Cell Line: "Prostate Cancer"	162088

The KM12L4, KM12C, and MDA-MB-231 cell lines are described in Example 1 above. The MCF7 cell line was derived from a pleural effusion of a breast adenocarcinoma and is non-metastatic. The MV-522 cell line is derived from a human lung carcinoma and is of high metastatic potential. The UCP-3 cell line is a low metastatic human lung carcinoma cell line; the MV-522 is a high metastatic variant of UCP-3. These cell lines are well-recognized in the art as models for the study of human breast and lung cancer (see, e.g., Chandrasekaran *et al.*, *Cancer Res.* (1979) 39:870 (MDA-MB-231 and MCF-7); Gastpar *et al.*, *J Med Chem* (1998) 41:4965 (MDA-MB-231 and

MCF-7); Ranson *et al.*, *Br J Cancer* (1998) 77:1586 (MDA-MB-231 and MCF-7); Kuang *et al.*, *Nucleic Acids Res* (1998) 26:1116 (MDA-MB-231 and MCF-7); Varki *et al.*, *Int J Cancer* (1987) 40:46 (UCP-3); Varki *et al.*, *Tumour Biol.* (1990) 11:327; (MV-522 and UCP-3); Varki *et al.*, *Anticancer Res.* (1990) 10:637; (MV-522); Kelner *et al.*, *Anticancer Res* (1995) 15:867 (MV-522); 5 and Zhang *et al.*, *Anticancer Drugs* (1997) 8:696 (MV522)). The samples of libraries 15-20 are derived from two different patients (UC#2, and UC#3). The bFGF-treated HMEC were prepared by incubation with bFGF at 10ng/ml for 2 hrs; the VEGF-treated HMEC were prepared by incubation with 20ng/ml VEGF for 2 hrs. Following incubation with the respective growth factor, the cells were washed and lysis buffer added for RNA preparation. The GRRpz and WOca cell lines were 10 provided by Dr. Donna M. Peehl, Department of Medicine, Stanford University School of Medicine. GRRpz was derived from normal prostate epithelium. The WOca cell line is a Gleason Grade 4 cell line.

Each of the libraries is composed of a collection of cDNA clones that in turn are 15 representative of the mRNAs expressed in the indicated mRNA source. In order to facilitate the analysis of the millions of sequences in each library, the sequences were assigned to clusters. The concept of "cluster of clones" is derived from a sorting/grouping of cDNA clones based on their hybridization pattern to a panel of roughly 300 7bp oligonucleotide probes (see Drmanac *et al.*, *Genomics* (1996) 37(1):29). Random cDNA clones from a tissue library are hybridized at moderate stringency to 300 7bp oligonucleotides. Each oligonucleotide has some measure of specific 20 hybridization to that specific clone. The combination of 300 of these measures of hybridization for 300 probes equals the "hybridization signature" for a specific clone. Clones with similar sequence will have similar hybridization signatures. By developing a sorting/grouping algorithm to analyze these signatures, groups of clones in a library can be identified and brought together computationally. These groups of clones are termed "clusters". Depending on the stringency of the 25 selection in the algorithm (similar to the stringency of hybridization in a classic library cDNA screening protocol), the "purity" of each cluster can be controlled. For example, artifacts of clustering may occur in computational clustering just as artifacts can occur in "wet-lab" screening of a cDNA library with 400 bp cDNA fragments, at even the highest stringency. The stringency used in the implementation of cluster herein provides groups of clones that are in general from the same 30 cDNA or closely related cDNAs. Closely related clones can be a result of different length clones of the same cDNA, closely related clones from highly related gene families, or splice variants of the same cDNA.

Differential expression for a selected cluster was assessed by first determining the number 35 of cDNA clones corresponding to the selected cluster in the first library (Clones in 1st), and the determining the number of cDNA clones corresponding to the selected cluster in the second library

(Clones in 2nd). Differential expression of the selected cluster in the first library relative to the second library is expressed as a "ratio" of percent expression between the two libraries. In general, the "ratio" is calculated by: 1) calculating the percent expression of the selected cluster in the first library by dividing the number of clones corresponding to a selected cluster in the first library by the 5 total number of clones analyzed from the first library; 2) calculating the percent expression of the selected cluster in the second library by dividing the number of clones corresponding to a selected cluster in a second library by the total number of clones analyzed from the second library; 3) dividing the calculated percent expression from the first library by the calculated percent expression from the second library. If the "number of clones" corresponding to a selected cluster in a library is 10 zero, the value is set at 1 to aid in calculation. The formula used in calculating the ratio takes into account the "depth" of each of the libraries being compared, *i.e.*, the total number of clones analyzed in each library.

In general, a polynucleotide is said to be significantly differentially expressed between two samples when the ratio value is greater than at least about 2, preferably greater than at least about 3, 15 more preferably greater than at least about 5, where the ratio value is calculated using the method described above. The significance of differential expression is determined using a z score test (Zar, Biostatistical Analysis, Prentice Hall, Inc., USA, "Differences between Proportions," pp 296-298 (1974).

20 Examples 5-12: Differential Expression of Polynucleotides of the Invention

A number of polynucleotide sequences have been identified that are differentially expressed between, for example, cells derived from high metastatic potential cancer tissue and low metastatic cancer cells, and between cells derived from high metastatic potential cancer tissue and normal tissue. Evaluation of the levels of expression of the genes corresponding to these sequences can be 25 valuable in diagnosis, prognosis, and/or treatment (*e.g.*, to facilitate rationale design of therapy, monitoring during and after therapy, *etc.*). Moreover, the genes corresponding to differentially expressed sequences described herein can be therapeutic targets due to their involvement in regulation (*e.g.*, inhibition or promotion) of development of, for example, the metastatic phenotype. For example, sequences that correspond to genes that are increased in expression in high metastatic 30 potential cells relative to normal or non-metastatic tumor cells may encode genes or regulatory sequences involved in processes such as angiogenesis, differentiation, cell replication, and metastasis.

Detection of the relative expression levels of differentially expressed polynucleotides described herein can provide valuable information to guide the clinician in the choice of therapy. 35 For example, a patient sample exhibiting an expression level of one or more of these polynucleotides

that corresponds to a gene that is increased in expression in metastatic or high metastatic potential cells may warrant more aggressive treatment for the patient. In contrast, detection of expression levels of a polynucleotide sequence that corresponds to expression levels associated with that of low metastatic potential cells may warrant a more positive prognosis than the gross pathology would suggest.

A number of polynucleotide sequences of the present invention are differentially expressed between human microvascular endothelial cells (HMEC) that have been treated with growth factors relative to untreated HMEC. Sequences that are differentially expressed between growth factor-treated HMEC and untreated HMEC can represent sequences encoding gene products involved in angiogenesis, metastasis (cell migration), and other development and oncogenic processes. For example, sequences that are more highly expressed in HMEC treated with growth factors (such as bFGF or VEGF) relative to untreated HMEC can serve as markers of cancer cells of higher metastatic potential. Detection of expression of these sequences in colon cancer tissue can be valuable in determining diagnostic, prognostic and/or treatment information associated with the prevention of achieving the malignant state in these tissues, and can be important in risk assessment for a patient. A patient sample displaying an increased level of one or more of these polynucleotides may thus warrant closer attention or more frequent screening procedures to catch the malignant state as early as possible.

The differential expression of the polynucleotides described herein can thus be used as, for example, diagnostic markers, prognostic markers, for risk assessment, patient treatment and the like. These polynucleotide sequences can also be used in combination with other known molecular and/or biochemical markers. The following examples provide relative expression levels of polynucleotides from specified cell lines and patient tissue samples.

25 **Example 5: High Metastatic Potential Breast Cancer Versus Low Metastatic Breast Cancer Cells**

The following tables summarize polynucleotides that represent genes that are differentially expressed between high metastatic potential and low metastatic potential breast cancer cells.

Table 5. High metastatic potential breast (lib3) > low metastatic potential (lib4) breast cancer cells

SEQ ID NO:	Lib3 Clones	Lib4 Clones	Lib3/Lib4
1213	40	0	39
1538	60	3	20
1466	14	0	14
1356	10	0	10
1383	10	1	10
1158	10	1	10
441	10	1	10
1338	10	0	10
1426	19	2	9

SEQ ID NO:	Lib3 Clones	Lib4 Clones	Lib3/Lib4
1547	9	1	9
1313	8	1	8
841	8	1	8
1534	8	0	8
1503	8	0	8
829	8	1	8
1408	8	0	8
1447	7	0	7
1389	7	0	7
356	7	0	7
1492	7	0	7
1543	22	3	7
799	7	0	7
1437	6	0	6
1251	6	0	6
972	18	3	6
1482	6	0	6
1299	6	0	6
109	24	4	6
1558	6	0	6
1355	6	0	6
1548	11	2	5
250	10	2	5
919	26	6	4
358	36	12	3
1525	75	28	3
1157	49	17	3

Table 6. Low metastatic potential breast (lib4) > high metastatic potential breast cancer cells (lib3)

SEQ ID NO:	Lib3 Clones	Lib4 Clones	Lib4/Lib3
248	0	58	59
726	1	23	24
14	1	19	19
699	0	14	14
763	1	14	14
20	1	13	13
79	1	13	13
715	0	10	10
991	0	8	8
1199	0	8	8
707	0	7	7
1128	4	26	7
891	0	6	6
1146	2	11	6
731	7	44	6
1518	3	15	5
340	3	13	4
949	4	13	3

SEQ ID NO:	Lib3 Clones	Lib4 Clones	Lib4/Lib3
1247	7	18	3
1185	497	1216	3

Example 6: High Metastatic Potential Lung Cancer Versus Low Metastatic Lung Cancer Cells

The following summarizes polynucleotides that represent genes differentially expressed between high metastatic potential lung cancer cells and low metastatic potential lung cancer cells:

5 Table 7. High metastatic potential lung (lib8) > low metastatic potential lung (lib9) lung cancer cells

SEQ ID NO:	Lib8 Clones	Lib9 Clones	Lib8/Lib9
150	31	0	43
651	43	2	30
1298	14	1	20
57	11	0	15
625	7	0	10
1322	7	1	10
36	7	0	10
621	18	3	8
215	6	1	8
561	19	4	7
247	5	0	7
199	5	0	7
998	5	0	7
502	5	0	7
1382	8	2	6
1181	17	4	6
1309	8	2	6
1157	15	4	5
1260	14	5	4
1185	710	266	4
1525	21	10	3

Table 8. Low metastatic potential lung (lib9) > high metastatic potential lung (lib8) cancer cells

SEQ ID NO:	Lib8 Clones	Lib9 Clones	Lib9/Lib8
924	1	13	9
822	1	13	9
728	1	12	9
341	1	12	9
1527	3	31	7
698	4	26	5
949	2	15	5
744	3	23	5
973	8	27	2

Example 7: High Metastatic Potential Colon Cancer Versus Low Metastatic Colon Cancer Cells

Tables 9 and 10 summarize polynucleotides that represent genes differentially expressed between high metastatic potential and low metastatic potential colon cancer cells:

5 **Table 9. High metastatic potential (lib1) > low metastatic potential (lib2) colon cancer cells**

SEQ ID NO:	Lib1 Clones	Lib2 Clones	Lib1/Lib2
248	67	2	31
87	12	0	11
698	11	0	10
57	13	3	4
924	24	10	2
1249	24	9	2

Table 10. Low metastatic potential (lib2) > high metastatic potential colon cancer (lib1) cells

SEQ ID NO:	Lib1 Clones	Lib2 Clones	Lib2/Lib1
1268	1	17	18
1114	0	15	16
1032	1	14	15
109	5	60	13
973	1	11	12
91	1	11	12
982	0	9	10
1267	3	28	10
93	1	8	9
1556	1	8	9
1251	0	8	9
1206	2	17	9
812	0	8	9
1254	0	7	8
1220	0	7	8
766	0	7	8
1156	0	7	8
1007	0	7	8
981	0	7	8
762	0	7	8
876	0	6	6
1234	2	11	6
1183	0	6	6
1044	2	12	6
785	0	6	6
1069	3	17	6
770	0	6	6
778	0	6	6
792	0	6	6
822	2	10	5
1258	7	23	4
1224	7	17	3

SEQ ID NO:	Lib1 Clones	Lib2 Clones	Lib2/Lib1
984	8	19	3
841	10	28	3
339	14	34	3
1213	11	29	3
1201	5	14	3
1192	22	48	2

Example 8: High Metastatic Potential Colon Cancer Patient Tissue Vs. Normal Patient Tissue

Tables 11 summarizes polynucleotides that represent genes differentially expressed between high metastatic potential colon cancer cells and normal colon cells of patient tissue.

5 **Table 11. High metastatic potential colon tissue (lib17) vs. normal colon tissue (lib15)**

SEQ ID NO:	Lib15 Clones	Lib17 Clones	Lib17/Lib15
1422	1	13	12
1132	1	10	9
730	1	9	8
1311	0	7	7
78	9	48	5
822	5	20	4
SEQ ID NO:	Lib15 Clones	Lib17 Clones	Lib15/Lib17
463	8	1	9

Example 9: High Tumor Potential Colon Tissue Vs. Metastasized Colon Cancer Tissue

The following table summarizes polynucleotides that represent genes differentially expressed between high tumor potential colon cancer cells and cells derived from high metastatic potential colon cancer cells of a patient.

10 **Table 12. High tumor potential colon tissue (lib16) vs. high metastatic colon tissue (lib17)**

SEQ ID NO:	Lib16 Clones	Lib17 Clones	Lib16/Lib17
1185	14	4	4
SEQ ID NO:	Lib16 Clones	Lib17 Clones	Lib17/Lib16
822	2	20	10

Example 10: High Tumor Potential Colon Cancer Patient Tissue Versus Normal Patient Tissue

Tables 13 and 14 summarize polynucleotides that represent genes differentially expressed between high metastatic potential colon cancer cells and normal colon cells in patient tissue:

15 **Table 13. Higher expression in tumor potential colon tissue (lib16) vs. normal colon tissue (lib15)**

SEQ ID NO:	Lib15 Clones	Lib16 Clones	Lib16/Lib15
1311	0	8	8
78	9	28	3

Table 14. Higher expression in normal colon tissue (lib15) vs. tumor potential colon tissue (lib16)

SEQ ID NO:	Lib15 Clones	Lib16 Clones	Lib15/Lib16
463	8	0	8
1099	12	3	4

Example 11: Growth Factor-Stimulated Human Microvascular Endothelial Cells (HMEC)**5 Relative to Untreated HMEC**

The following tables summarize polynucleotides that represent genes differentially expressed between growth factor-treated and untreated HMEC.

Table 15. Higher expression in bFGF treated HMEC (lib13) vs. untreated HMEC (lib12)

SEQ ID NO:	Lib12 Clones	Lib13 Clones	Lib13/Lib12
1520	9	23	3
1538	17	35	2

10 Table 16. Higher expression in VEGF treated HMEC (lib14) vs. untreated HMEC (lib12)

SEQ ID NO:	Lib12 Clones	Lib14 Clones	Lib14/Lib12
1154	2	12	6
1226	2	10	5
1538	17	38	2

Example 12: Polynucleotides Differentially Expressed in Human Prostate Cancer Cells Relative to Normal Human Prostate Cells

The following tables summarize identified polynucleotides that represent genes differentially expressed between prostate cancer cells and normal prostate cells:

15 Table 17. Higher expression in normal prostate cells (lib21) relative to prostate cancer cells (lib22)

SEQ ID NO:	Lib21 Clones	Lib22 Clones	Lib21/Lib22
1525	6	0	6
248	116	51	2
1203	22	9	2

Table 18 Higher expression in prostate cancer cells (lib22) relative to normal prostate cells (lib21)

SEQ ID NO:	Lib21 Clones	Lib22 Clones	Lib22/Lib21
1213	0	34	35
340	1	12	12
699	0	11	11

20 Example 13: Differential Expression Across Multiple Libraries

A number of polynucleotide sequences have been identified that represent genes that are differentially expressed across multiple libraries. Expression of these sequences in a tissue or any

origin can be valuable in determining diagnostic, prognostic and/or treatment information associated with the prevention of achieving the malignant state in these tissues, and can be important in risk assessment for a patient. These polynucleotides can also serve as non-tissue specific markers of, for example, risk of metastasis of a tumor. Table 19 summarizes this data.

5

Table 19. Genes Differentially Expressed Across Multiple Library Comparisons

SEQ ID NO:	Cell or Tissue Sample and Cancer State Compared	Ratio
57	High Met Lung (lib8) > Low Met Lung (lib9)	15
57	High Met Colon (lib1) > Low Met Colon (lib2)	4
78	High Met Colon Tissue (lib17) > Normal Colon Tissue (lib15)	5
78	Normal Colon Tumor Tissue (lib16) > Normal Colon Tissue (lib15)	3
109	High Met Breast (lib3) > Low Met Breast (lib4)	6
109	Low Met Colon (lib2) > High Met Colon (lib1)	13
248	High Met Colon (lib1) > Low Met Colon (lib2)	31
248	Normal Prostate (lib21) > Prostate Cancer (lib22)	2
248	Low Met Breast (lib4) > High Met Breast (lib3)	59
340	Prostate Cancer (lib22) > Normal Prostate (lib21)	12
340	Low Met Breast (lib4) > High Met Breast (lib3)	4
463	Normal Colon Tissue (lib15) > High Met Colon Tissue (lib17)	9
463	Normal Colon Tissue (lib15) > Normal Colon Tumor Tissue (lib16)	8
698	High Met Colon (lib1) > Low Met Colon (lib2)	10
698	Low Met Lung (lib9) > High Met Lung (lib8)	5
699	Low Met Breast (lib4) > High Met Breast (lib3)	14
699	Prostate Cancer (lib22) > Normal Prostate (lib21)	11
822	High Met Colon Tissue (lib17) > Normal Colon Tumor Tissue (lib16)	10
822	Low Met Lung (lib9) > High Met Lung (lib8)	9
822	Low Met Colon (lib2) > High Met Colon (lib1)	5
822	High Met Colon Tissue (lib17) > Normal Colon Tissue (lib15)	4
841	High Met Breast (lib3) > Low Met Breast (lib4)	8
841	Low Met Colon (lib2) > High Met Colon (lib1)	3
924	High Met Colon (lib1) > Low Met Colon (lib2)	2
924	Low Met Lung (lib9) > High Met Lung (lib8)	9
949	Low Met Lung (lib9) > High Met Lung (lib8)	5
949	Low Met Breast (lib4) > High Met Breast (lib3)	3
973	Low Met Colon (lib2) > High Met Colon (lib1)	12
973	Low Met Lung (lib9) > High Met Lung (lib8)	2
1157	High Met Lung (lib8) > Low Met Lung (lib9)	5
1157	High Met Breast (lib3) > Low Met Breast (lib4)	3
1185	Normal Colon Tumor Tissue (lib16) > High Met Colon Tissue (lib17)	4
1185	High Met Lung (lib8) > Low Met Lung (lib9)	4
1185	Low Met Breast (lib4) > High Met Breast (lib3)	3
1213	High Met Breast (lib3) > Low Met Breast (lib4)	39
1213	Prostate Cancer (lib22) > Normal Prostate (lib21)	35
1213	Low Met Colon (lib2) > High Met Colon (lib1)	3
1251	High Met Breast (lib3) > Low Met Breast (lib4)	6
1251	Low Met Colon (lib2) > High Met Colon (lib1)	9
1311	Normal Colon Tumor Tissue (lib16) > Normal Colon Tissue (lib15)	8

SEQ ID NO:	Cell or Tissue Sample and Cancer State Compared	Ratio
1311	High Met Colon Tissue (lib17) > Normal Colon Tissue (lib15)	7
1525	Normal Prostate (lib21) > Prostate Cancer (lib22)	6
1525	High Met Lung (lib8) > Low Met Lung (lib9)	3
1525	High Met Breast (lib3) > Low Met Breast (lib4)	3
1538	High Met Breast (lib3) > Low Met Breast (lib4)	20
1538	HMEC-VEGF (lib14) > HMEC (lib12)	2
1538	HMEC-bFGF (lib13) > HMEC (lib12)	2

Key for Table 19: High Met = high metastatic potential; Low Met = low metastatic potential; met = metastasized; tumor = non-metastasized tumor; HMEC = human microvascular endothelial cell; bFGF = bFGF treated; VEGF = VEGF treated.

5 Example 14: Identification of Contiguous Sequences Having a Polynucleotide of the Invention

The novel polynucleotides were used to screen publicly available and proprietary databases to determine if any of the polynucleotides of SEQ ID NOS:2611-2707 would facilitate identification of a contiguous sequence, e.g., the polynucleotides would provide sequence that would result in 5' extension of another DNA sequence, resulting in production of a longer contiguous sequence

10 composed of the provided polynucleotide and the other DNA sequence(s). Contiging was performed using the Gelmerge application (default settings) of GCG from the Univ. of Wisconsin.

15 Using these parameters, 97 contiged sequences were generated. These contiged sequences are provided as SEQ ID NOS:2611-2707 (see Table 1C). Table 1C provides the SEQ ID NO of the contig sequence, the name of the sequence used to create the contig, and the accession number of the publicly available tentative human consensus (THC) sequence used with the sequence of the corresponding sequence name to provide the contig. The sequence name of Table 1C can be correlated with the SEQ ID NO: of the polynucleotide of the invention using Tables 1A and 1B.

20 The contiged sequences (SEQ ID NOS:2611-2707) thus represent longer sequences that encompass a polynucleotide sequence of the invention. The contiged sequences were then translated in all three reading frames to determine the best alignment with individual sequences using the BLAST programs as described above. The sequences were masked using the XBLAST program for masking low complexity as described above in Example 1. Several of the contiged sequences were found to encode polypeptides having characteristics of a polypeptide belonging to a known protein families (and thus represent new members of these protein families) and/or comprising a known functional domain (Table 3B, inserted prior to claims). Thus the invention encompasses fragments, fusions, and variants of such polynucleotides that retain biological activity associated with the protein family and/or functional domain identified herein.

25 Descriptions of the profiles for the indicated protein families and functional domains are provided in Example 3 above. A description of the profile for PR55 is provided below.

Protein Phosphatase 2A Regulatory Subunit PR55 (PR55). Several of the contigs correspond to a sequence encoding a protein comprising a protein phosphatase 2A (PP2A) regulatory subunit PR55. PP2A is a serine/threonine phosphatase involved in many aspects of cellular function including the regulation of metabolic enzymes and proteins involved in signal transduction. PP2A is a trimeric enzyme comprising a core composed of a catalytic subunit associated with a 65 Kd regulatory subunit (PR65, also called subunit A). This complex associates with a third variable subunit (subunit B), which confers distinct properties to the holoenzyme (Mayer-Jaekel et al. *Trends Cell Biol.* (1994) 4:287-291). One of the forms of the variable subunit is a 55 Kd protein (PR55) which is highly conserved in mammals and may facilitate substrate recognition or targeting the enzyme complex to the appropriate subcellular compartment. The PR55 subunit comprises two conserved sequences of 15 residues; one located in the N-terminal region, the other in the center of the protein. The consensus patterns are: E-F-D-Y-L-K-S-L-E-I-E-E-K-I-N; and N-[AG]-H-[TA]-Y-H-I-N-S-I-S-[LIVM]-N-S-D.

Those skilled in the art will recognize, or be able to ascertain, using not more than routine experimentation, many equivalents to the specific embodiments of the invention described herein. Such specific embodiments and equivalents are intended to be encompassed by the following claims.

All publications and patent applications cited in this specification are herein incorporated by reference as if each individual publication or patent application were specifically and individually indicated to be incorporated by reference. The citation of any publication is for its disclosure prior to the filing date and should not be construed as an admission that the present invention is not entitled to antedate such publication by virtue of prior invention.

Although the foregoing invention has been described in some detail by way of illustration and example for purposes of clarity of understanding, it is readily apparent to those of ordinary skill in the art in light of the teachings of this invention that certain changes and modifications may be made thereto without departing from the spirit or scope of the appended claims.

Deposit Information. The following materials were deposited with the American Type Culture Collection (CMCC = Chiron Master Culture Collection).

30 **Table 20. Cell Lines Deposited with ATCC**

Cell Line	Deposit Date	ATCC Accession No.	CMCC Accession No.
KM12L4-A	March 19, 1998	CRL-12496	11606
Km12C	May 15, 1998	CRL-12533	11611
MDA-MB-231	May 15, 1998	CRL-12532	10583
MCF-7	October 9, 1998	CRL-12584	10377

In addition, pools of selected clones, as well as libraries containing specific clones, were assigned an "ES" number (internal reference) and deposited with the ATCC. Table 21 below provides the ATCC Accession Nos. of the ES deposits, all of which were deposited on or before May 13, 1999. The names of the clones contained within each of these deposits are provided in the 5 tables numbered 22 and greater (inserted before the claims).

Table 21: Pools of Clones and Libraries Deposited with ATCC on or before May 14, 1999

ES #	ATCC Accession #	ES #	ATCC Accession #	ES #	ATCC Accession #
34		41		48	
35		42		49	
36		43		50	
37		44		51	
38		45		52	
39		46		53	
40		47		54	

The deposits described herein are provided merely as convenience to those of skill in the art, and is not an admission that a deposit is required under 35 U.S.C. §112. The sequence of the 10 polynucleotides contained within the deposited material, as well as the amino acid sequence of the polypeptides encoded thereby, are incorporated herein by reference and are controlling in the event of any conflict with the written description of sequences herein. A license may be required to make, use, or sell the deposited material, and no such license is granted hereby.

Retrieval of Individual Clones from Deposit of Pooled Clones. Where the ATCC deposit is composed of a pool of cDNA clones or a library of cDNA clones, the deposit was prepared by first 15 transfecting each of the clones into separate bacterial cells. The clones in the pool or library were then deposited as a pool of equal mixtures in the composite deposit. Particular clones can be obtained from the composite deposit using methods well known in the art. For example, a bacterial cell containing a particular clone can be identified by isolating single colonies, and identifying colonies containing the specific clone through standard colony hybridization techniques, using an oligonucleotide probe or 20 probes designed to specifically hybridize to a sequence of the clone insert (e.g., a probe based upon unmasked sequence of the encoded polynucleotide having the indicated SEQ ID NO). The probe should be designed to have a T_m of approximately 80°C (assuming 2°C for each A or T and 4°C for each G or C). Positive colonies can then be picked, grown in culture, and the recombinant clone isolated. Alternatively, probes designed in this manner can be used to PCR to isolate a nucleic acid molecule 25 from the pooled clones according to methods well known in the art, e.g., by purifying the cDNA from the deposited culture pool, and using the probes in PCR reactions to produce an amplified product having the corresponding desired polynucleotide sequence.

We Claim:

1. A library of polynucleotides, the library comprising the sequence information of at least one of SEQ ID NOS:1-2702.

5

2. The library of claim 1, wherein the library is provided on a nucleic acid array.

3. The library of claim 1, wherein the library is provided in a computer-readable format.

10 4. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in a cancer cell of high metastatic potential relative to a control cell, wherein the control cell is a normal cell or a cell of low metastatic potential, and wherein the sequence is selected from the group consisting of SEQ ID NOS:1213, 1538, 1466, 1356, 1383, 1158, 441, 1338, 1426, 1547, 1313, 841, 1534, 1503, 829, 1408, 1447, 1389, 356, 1492, 1543, 799, 1437, 1251, 972, 1482, 1299, 109, 1558, 1355, 1548, 250, 919, 358, 1525, 1157, 150, 651, 1298, 15 57, 625, 1322, 36, 621, 215, 561, 247, 199, 998, 502, 1382, 1181, 1309, 1157, 1260, 1185, 1525, 248, 87, 698, 57, 924, 1249.

20 5. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in a cancer cell of low metastatic potential relative to a control cell, wherein the control cell is a normal cell or a cell of high metastatic potential, and wherein the sequence is selected from the group consisting of SEQ ID NOS:248, 726, 14, 699, 763, 20, 79, 715, 991, 1199, 707, 1128, 891, 1146, 731, 1518, 340, 949, 1247, 1185, 924, 822, 728, 341, 1527, 698, 949, 744, 973, 1268, 1114, 1032, 109, 973, 91, 982, 1267, 93, 1556, 1251, 1206, 812, 1254, 1220, 25 766, 1156, 1007, 981, 762, 876, 1234, 1183, 1044, 785, 1069, 770, 778, 792, 822, 1258, 1224, 984, 841, 339, 1213, 1201, 1192.

30 6. An isolated polynucleotide comprising a nucleotide sequence having at least 90% sequence identity to an identifying sequence of SEQ ID NOS:1-2707 or a degenerate variant or fragment thereof.

7. A recombinant host cell containing the polynucleotide of claim 6.

8. An isolated polypeptide encoded by the polynucleotide of claim 6.

35

9. An antibody that specifically binds a polypeptide of claim 8.

10. A vector comprising the polynucleotide of claim 6.

11. A polynucleotide comprising the nucleotide sequence of an insert contained in a clone deposited as ATCC accession number xx, xx, xx, xx, xx, xx, xx, xx, or xx.

5

12. A method of detecting differentially expressed genes correlated with a cancerous state of a mammalian cell, the method comprising the step of:

detecting at least one differentially expressed gene product in a test sample derived from a cell suspected of being cancerous, where the gene product is encoded by a gene corresponding to a sequence of at least one of SEQ ID NOS: 1213, 1538, 1466, 1356, 1383, 1158, 441, 1338, 1426, 1547, 1313, 841, 1534, 1503, 829, 1408, 1447, 1389, 356, 1492, 1543, 799, 1437, 1251, 972, 1482, 1299, 109, 1558, 1355, 1548, 250, 919, 358, 1525, 1157, 150, 651, 1298, 57, 625, 1322, 36, 621, 215, 561, 247, 199, 998, 502, 1382, 1181, 1309, 1157, 1260, 1185, 1525, 248, 87, 698, 57, 924, 1249, 248, 726, 14, 699, 763, 20, 79, 715, 991, 1199, 707, 1128, 891, 1146, 731, 1518, 340, 949, 1247, 1185, 924, 822, 728, 341, 1527, 698, 949, 744, 973, 1268, 1114, 1032, 109, 973, 91, 982, 1267, 93, 1556, 1251, 1206, 812, 1254, 1220, 766, 1156, 1007, 981, 762, 876, 1234, 1183, 1044, 785, 1069, 770, 778, 792, 822, 1258, 1224, 984, 841, 339, 1213, 1201, 1192

wherein detection of the differentially expressed gene product is correlated with a cancerous state of the cell from which the test sample was derived.